

# Provable Pathways: Learning Multiple Tasks over Multiple Paths

Yingcong Li\*

Samet Oymak\*†

\* University of California, Riverside

† University of Michigan, Ann Arbor

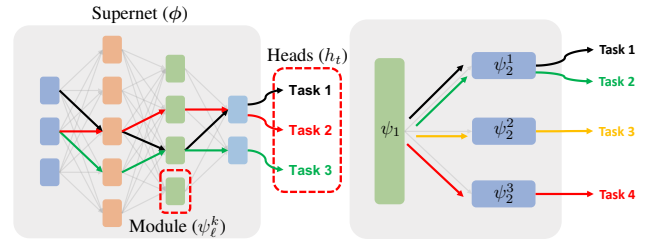
## Abstract

Constructing useful representations across a large number of tasks is a key requirement for sample-efficient intelligent systems. A traditional idea in multitask learning (MTL) is building a shared representation across tasks which can then be adapted to new tasks by tuning last layers. A desirable refinement of using a shared one-fits-all representation is to construct task-specific representations. To this end, recent PathNet/muNet architectures represent individual tasks as pathways within a larger supernet. The subnetworks induced by pathways can be viewed as task-specific representations that are composition of modules within supernet’s computation graph. This work explores the pathways proposal from the lens of statistical learning: We first develop novel generalization bounds for empirical risk minimization problems learning multiple tasks over multiple paths (Multipath MTL). In conjunction, we formalize the benefits of resulting multipath representation when adapting to new downstream tasks. Our bounds are expressed in terms of Gaussian complexity, lead to tangible guarantees for the class of linear representations, and provide novel insights into the quality and benefits of a multipath representation. When computation graph is a tree, Multipath MTL hierarchically clusters the tasks and builds cluster-specific representations. We provide further discussion and experiments for hierarchical MTL and rigorously identify the conditions under which Multipath MTL is provably superior to traditional MTL approaches with shallow supernet.

## 1 Introduction

Multitask learning (MTL) promises to deliver significant accuracy improvements by leveraging similarities across many tasks through shared representations. The potential of MTL has been recognized since 1990s (Caruana 1997) however its impact has grown over time thanks to more recent machine learning applications arising in computer vision and NLP that involve large datasets with thousands of classes/tasks. Representation learning techniques (e.g. MTL and self-supervision) are also central to the success of deep learning as large pre-trained models enable data-efficient learning for downstream transfer learning tasks (Deng et al. 2009; Brown et al. 2020).

As we move from tens of tasks trained with small models to thousands of tasks trained with large models, new



(a) General computation graph

(b) Hierarchical MTL

Figure 1: In Multipath MTL, each task selects a pathway within a supernet graph. The composition of the modules along the pathway forms the task-specific representation. Fig. 1a depicts a general supernet graph (highlighted in gray block), and the pathways for different tasks are shown in colored arrows. Fig. 1b is a special instance where related tasks are hierarchically clustered: For instance, Tasks 1 and 2 are assigned the same representation  $\psi_2^1 \circ \psi_1$ .

statistical and computational challenges arise: First, not all tasks will be closely related to each other, for instance, tasks might admit a natural clustering into groups. This is also connected to heterogeneity challenge in federated learning where clients have distinct distributions and benefit from personalization. To address this challenge, rather than a single task-agnostic representation, it might be preferable to use a task-specific representation. Secondly, pretrained language and vision models achieve better accuracy with larger sizes which creates computational challenges as they push towards trillion parameters. This motivated new architectural proposals such as Pathways/PathNet (Fernando et al. 2017; Dean 2021; Gesmundo and Dean 2022b) where tasks can be computed over compute-efficient subnetworks. At a high-level, each subnetwork is created by a composition of modules within a larger supernet which induces a pathway as depicted in Figure 1. Inspired from these challenges, we ask

**Q:** What are the statistical benefits of learning task-specific representations along supernet pathways?

Our primary contribution is formalizing the Multipath MTL problem depicted in Figure 1 and developing associated statistical learning guarantees that shed light on its benefits. Our formulation captures important aspects of the problem including learning compositional MTL representations, multilayer

\*Emails: {yli692@,oymak@ece.}ucr.edu

nature of supernet, assigning optimal pathways to individual tasks, and transferring learned representations to novel downstream tasks. Our specific contributions are as follows.

- Suppose we have  $N$  samples per task and  $T$  tasks in total. Denote the hypothesis sets for multipath representation by  $\Phi$ , task specific heads by  $\mathcal{H}$  and potential pathway choices by  $\mathcal{A}$ . Our main result bounds the task-averaged risk of MTL as

$$\sqrt{\frac{\text{DoF}(\Phi_{\text{used}})}{NT}} + \sqrt{\frac{\text{DoF}(\mathcal{H}) + \text{DoF}(\mathcal{A})}{N}}. \quad (1)$$

Here,  $\text{DoF}(\cdot)$  returns the *degrees of freedom* of a hypothesis set (i.e. number of parameters). More generally, Theorem 1 states our guarantees in terms of Gaussian complexity.  $\Phi_{\text{used}} \subseteq \Phi$  is the supernet spanned by the pathways of the empirical solution and  $1/NT$  dependence implies that cost of representation learning is shared across tasks. We also show a *no-harm* result (Lemma 1): If the supernet is sufficiently expressive to achieve zero empirical risk, then, the excess risk of individual tasks will not be harmed by the other tasks. Theorem 2 develops guarantees for transferring the resulting MTL representation to a new task in terms of representation bias of the empirical MTL supernet.

- When the supernet has a single module, the problem boils down to (vanilla) MTL with single shared representation and our bounds recover the results by (Maurer, Pontil, and Romera-Paredes 2016; Tripuraneni, Jin, and Jordan 2021). When the supernet graph is hierarchical (as in Figure 1b), our bounds provide insights for the benefits of clustering tasks into similar groups and superiority of multilayer Multipath MTL over using single-layer shallow supernets (Section 5).
- We develop stronger results for linear representations over a supernet and obtain novel MTL and transfer learning bounds (Sec. 4 and Theorem 4). These are accomplished by developing new task-diversity criteria to account for the task-specific (thus heterogeneous) nature of multipath representations. Numerical experiments support our theory and verify the benefits of multipath representations. Finally, we also highlight multiple future directions.

## 2 Setup and Problem Formulations

**Notation.** Let  $\|\cdot\|$  denote the  $\ell_2$ -norm of a vector and operator norm of a matrix.  $|\cdot|$  denotes the absolute value for scalars and cardinality for discrete sets. We use  $[K]$  to denote the set  $\{1, 2, \dots, K\}$  and  $\lesssim, \gtrsim$  for inequalities that hold up to constant/logarithmic factors.  $\mathcal{Q}^K$  denotes  $K$ -times Cartesian product of a set  $\mathcal{Q}$  with itself.  $\circ$  denotes functional composition, i.e.,  $f \circ g(x) = f(g(x))$ .

**Setup.** Suppose we have  $T$  tasks each following data distribution  $\{\mathcal{D}_t\}_{t=1}^T$ . During MTL phase, we are given  $T$  training datasets  $\{\mathcal{S}_t\}_{t=1}^T$  each drawn i.i.d. from its corresponding distribution  $\mathcal{D}_t$ . Let  $\mathcal{S}_t = \{(\mathbf{x}_{ti}, y_{ti})\}_{i=1}^N$ , where  $(\mathbf{x}_{ti}, y_{ti}) \in (\mathcal{X}, \mathbb{R})$  is an input-label pair and  $\mathcal{X}$  is the input space, and  $|\mathcal{S}_t| = N$  is the number of samples per task. We assume the same  $N$  for all tasks for cleaner exposition. Define the union of the datasets by  $\mathcal{S}_{\text{all}} = \bigcup_{t=1}^T \mathcal{S}_t$  (with  $|\mathcal{S}_{\text{all}}| = NT$ ), and the set of distributions by  $\mathcal{D} = \{\mathcal{D}_t\}_{t=1}^T$ .

Following the setting of related works (Tripuraneni, Jin, and Jordan 2021), we will consider two problems: **(1) MTL**

**problem** will use these  $T$  datasets to learn a supernet and establish guarantees for representation learning. **(2) Transfer learning problem** will use the resulting representation for a downstream task in a sample efficient fashion.

**Problem (1): Multipath Multitask Learning (M<sup>2</sup>TL).** We consider a supernet with  $L$  layers where layer  $\ell$  has  $K_\ell$  modules for  $\ell \in [L]$ . As depicted in Figure 1, each task will compose a task-specific representation by choosing one module from each layer. We refer to each sequence of  $L$  modules as a *pathway*. Let  $\mathcal{A} = [K_1] \times \dots \times [K_L]$  be the set of all pathway choices obeying  $|\mathcal{A}| = \prod_{\ell=1}^L K_\ell$ . Let  $\alpha_t \in \mathcal{A}$  denote the pathway associated with task  $t \in [T]$  where  $\alpha_t[\ell] \in [K_\ell]$  denotes the selected module index from layer  $\ell$ . We remark that results can be extended to more general pathway sets as discussed in Section 3.1.

As depicted in Figure 1, let  $\Psi_\ell$  be the hypothesis set of modules in  $\ell_{\text{th}}$  layer and  $\psi_\ell^k \in \Psi_\ell$  denote the  $k_{\text{th}}$  module function in the  $\ell_{\text{th}}$  layer, referred to as  $(\ell, k)$ 'th module. Let  $h_t \in \mathcal{H}$  be the prediction head of task  $t$  where all tasks use the same hypothesis set  $\mathcal{H}$  for prediction. Let us denote the combined hypothesis

$$\begin{aligned} \mathbf{h} &= [h_1, \dots, h_T] \in \mathcal{H}^T, \\ \boldsymbol{\alpha} &= [\alpha_1, \dots, \alpha_T] \in \mathcal{A}^T, \\ \boldsymbol{\psi}_\ell &= [\psi_\ell^1, \dots, \psi_\ell^{K_\ell}] \in \Psi_\ell^{K_\ell}, \forall \ell \in [L], \\ \boldsymbol{\phi} &:= [\boldsymbol{\psi}_1, \dots, \boldsymbol{\psi}_L] \in \Phi \end{aligned}$$

where  $\Phi = \Psi_1^{K_1} \times \dots \times \Psi_L^{K_L}$  is the supernet hypothesis class containing all modules/layers. Given a supernet  $\boldsymbol{\phi} \in \Phi$  and pathway  $\alpha$ ,  $\boldsymbol{\phi}_\alpha = \psi_L^\alpha \circ \dots \circ \psi_1^\alpha$  denotes the representation induced by pathway  $\alpha$  where we use the convention  $\psi_\ell^\alpha := \psi_\ell^{\alpha[\ell]}$ . Hence,  $\boldsymbol{\phi}_{\alpha_t}$  is the representation of task  $t$ . We would like to solve for supernet weights  $\boldsymbol{\phi}$ , pathways  $\boldsymbol{\alpha}$ , and heads  $\mathbf{h}$ . Thus, given a loss function  $\ell(\hat{y}, y)$ , Multipath MTL (M<sup>2</sup>TL) solves the following empirical risk minimization problem over  $\mathcal{S}_{\text{all}}$  to optimize the combined hypothesis  $\mathbf{f} = (\mathbf{h}, \boldsymbol{\alpha}, \boldsymbol{\phi})$ :

$$\hat{\mathbf{f}} = \arg \min_{\mathbf{f} \in \mathcal{F}} \widehat{\mathcal{L}}_{\mathcal{S}_{\text{all}}}(\mathbf{f}) := \frac{1}{T} \sum_{t=1}^T \widehat{\mathcal{L}}_t(h_t \circ \boldsymbol{\phi}_{\alpha_t}) \quad (\text{M}^2\text{TL})$$

$$\text{where } \widehat{\mathcal{L}}_t(f) = \frac{1}{N} \sum_{i=1}^N \ell(f(\mathbf{x}_{ti}), y_{ti})$$

$$\mathcal{F} := \mathcal{H}^T \times \mathcal{A}^T \times \Phi.$$

Here  $\widehat{\mathcal{L}}_t$  and  $\widehat{\mathcal{L}}_{\mathcal{S}_{\text{all}}}$  are task-conditional and task-averaged empirical risks. We are primarily interested in controlling the task-averaged test risk  $\mathcal{L}_{\overline{\mathcal{D}}}(\mathbf{f}) = \mathbb{E}[\widehat{\mathcal{L}}_{\mathcal{S}_{\text{all}}}(\mathbf{f})]$ . Let  $\mathcal{L}_{\overline{\mathcal{D}}}^* := \min_{\mathbf{f} \in \mathcal{F}} \mathcal{L}_{\overline{\mathcal{D}}}(\mathbf{f})$ , then the *excess MTL risk* is defined as

$$\mathcal{R}_{\text{M}^2\text{TL}}(\hat{\mathbf{f}}) = \mathcal{L}_{\overline{\mathcal{D}}}(\hat{\mathbf{f}}) - \mathcal{L}_{\overline{\mathcal{D}}}^*. \quad (2)$$

**Problem (2): Transfer Learning with Optimal Pathway (TLOP).** Suppose we have a novel target task with i.i.d. training dataset  $\mathcal{S}_{\mathcal{T}} = \{(\mathbf{x}_i, y_i)\}_{i=1}^M$  with  $M$  samples drawn from distribution  $\mathcal{D}_{\mathcal{T}}$ . Given a pretrained supernet  $\boldsymbol{\phi}$  (e.g., following (M<sup>2</sup>TL)), we can search for a pathway  $\alpha$  so that  $\boldsymbol{\phi}_\alpha$  becomes a suitable representation for  $\mathcal{D}_{\mathcal{T}}$ . Thus, for this new

task, we only need to optimize the path  $\alpha \in \mathcal{A}$  and the prediction head  $h \in \mathcal{H}_{\mathcal{T}}$  while reusing weights of  $\phi$ . This leads to the following problem:

$$\hat{f}_{\phi} = \arg \min_{h \in \mathcal{H}_{\mathcal{T}}, \alpha \in \mathcal{A}} \widehat{\mathcal{L}}_{\mathcal{T}}(f) \quad \text{where } f = h \circ \phi_{\alpha} \quad (\text{TLOP})$$

$$\text{and } \widehat{\mathcal{L}}_{\mathcal{T}}(f) = \frac{1}{M} \sum_{i=1}^M \ell(f(\mathbf{x}_i), y_i).$$

Here,  $\hat{f}_{\phi}$  reflects the fact that solution depends on the suitability of pretrained supernet  $\phi$ . Let  $f_{\phi}^*$  be a population minima of (TLOP) given supernet  $\phi$  (as  $M \rightarrow \infty$ ) and define the population risk  $\mathcal{L}_{\mathcal{T}}(f) = \mathbb{E}[\widehat{\mathcal{L}}_{\mathcal{T}}(f)]$ . (TLOP) will be evaluated against the hindsight knowledge of optimal supernet for target: Define the optimal target risk  $\mathcal{L}_{\mathcal{T}}^* := \min_{h \in \mathcal{H}_{\mathcal{T}}, \phi \in \Phi} \mathcal{L}_{\mathcal{T}}(h \circ \phi_{\alpha})$  which optimizes  $h, \phi$  for the target task along the fixed pathway  $\alpha = [1, \dots, 1]$ . Here we can fix  $\alpha$  since all pathways result in the same search space. We define the *excess transfer learning risk* to be

$$\begin{aligned} \mathcal{R}_{\text{TLOP}}(\hat{f}_{\phi}) &= \mathcal{L}_{\mathcal{T}}(\hat{f}_{\phi}) - \mathcal{L}_{\mathcal{T}}^* \\ &= \underbrace{\mathcal{L}_{\mathcal{T}}(\hat{f}_{\phi}) - \mathcal{L}_{\mathcal{T}}(f_{\phi}^*)}_{\text{variance}} + \underbrace{\mathcal{L}_{\mathcal{T}}(f_{\phi}^*) - \mathcal{L}_{\mathcal{T}}^*}_{\text{supernet bias}}. \end{aligned} \quad (3)$$

The final line decomposes the overall risk into a *variance* term and *supernet bias*. The former arises from the fact that we solve the problem with finite training samples. This term will vanish as  $M \rightarrow \infty$ . The latter term quantifies the bias induced by the fact that (TLOP) uses the representation  $\phi$  rather than the optimal representation. Finally, while supernet  $\phi$  in (TLOP) is arbitrary, for end-to-end guarantees we will set it to the solution  $\hat{\phi}$  of (M<sup>2</sup>TL). In this scenario, we will refer to  $\{\mathcal{D}_t\}_{t=1}^T$  as source tasks.

### 3 Main Results

We are ready to present our results that establish generalization guarantees for multitask and transfer learning problems over supernet pathways. Our results will be stated in terms of Gaussian complexity which is introduced below.

**Definition 1 (Gaussian Complexity)** Let  $\mathcal{Q}$  be a set of hypotheses that map  $\mathcal{Z}$  to  $\mathbb{R}^r$ . Let  $(\mathbf{g}_i)_{i=1}^n$  ( $\mathbf{g}_i \in \mathbb{R}^r$ ) be  $n$  independent vectors each distributed as  $\mathcal{N}(\mathbf{0}, \mathbf{I}_r)$  and let  $\mathbf{Z} = (\mathbf{z}_i)_{i=1}^n \in \mathcal{Z}^n$  be a dataset of input features. Then, the empirical Gaussian complexity is defined as

$$\widehat{\mathcal{G}}_{\mathbf{Z}}(\mathcal{Q}) = \mathbb{E}_{\mathbf{g}_i} \left[ \sup_{q \in \mathcal{Q}} \frac{1}{n} \sum_{i=1}^n \mathbf{g}_i^{\top} q(\mathbf{z}_i) \right].$$

The worst-case Gaussian complexity is obtained by considering the supremum over  $\mathbf{Z} \in \mathcal{Z}^n$  as follows

$$\widetilde{\mathcal{G}}_n^{\mathcal{Z}}(\mathcal{Q}) = \sup_{\mathbf{Z} \in \mathcal{Z}^n} [\widehat{\mathcal{G}}_{\mathbf{Z}}(\mathcal{Q})].$$

For cleaner notation, we drop the superscript  $\mathcal{Z}$  from the worst-case Gaussian complexity (using  $\widetilde{\mathcal{G}}_n(\mathcal{Q})$ ) as its input space will be clear from context. When  $\mathbf{Z} = (\mathbf{z}_i)_{i=1}^n$  are

drawn i.i.d. from  $\mathcal{D}$ , the (usual) Gaussian complexity is defined by  $\mathcal{G}_n(\mathcal{Q}) = \mathbb{E}_{\mathbf{Z} \sim \mathcal{D}^n} [\widehat{\mathcal{G}}_{\mathbf{Z}}(\mathcal{Q})]$ . Note that, we always have  $\mathcal{G}_n(\mathcal{Q}) \leq \widetilde{\mathcal{G}}_n(\mathcal{Q})$  assuming  $\mathcal{D}$  is supported on  $\mathcal{Z}$ . In our setting, keeping track of distributions along exponentially many pathways proves challenging, and we opt to use  $\widetilde{\mathcal{G}}_n(\mathcal{Q})$  which leads to clean upper bounds. The supplementary material also derives tighter but more convoluted bounds in terms of empirical complexity. Finally, it is well-known that Gaussian/Rademacher complexities scale as  $\sqrt{\text{comp}(\mathcal{Q})/n}$  where  $\text{comp}(\mathcal{Q})$  is a set complexity such as VC-dimension, which links to our informal statement (1).

We will first present our generalization bounds for the Multipath MTL problem using empirical process theory arguments. Our bounds will lead to meaningful guarantees for specific MTL settings, including vanilla MTL where all tasks share a single representation, as well as hierarchical MTL depicted in Fig. 1b. We will next derive transfer learning guarantees in terms of supernet bias, which quantifies the performance difference of a supernet from its optimum for a target. To state our results, we introduce two standard assumptions.

**Assumption 1** Elements of hypothesis sets  $\mathcal{H}$  and  $(\Psi_{\ell})_{\ell=1}^L$  are  $\Gamma$ -Lipschitz functions with respect to Euclidean norm.

**Assumption 2** Loss function  $\ell(\cdot, y) : \mathbb{R} \times \mathbb{R} \rightarrow [0, 1]$  is  $\Gamma$ -Lipschitz with respect to Euclidean norm.

#### 3.1 Results for Multipath Multitask Learning

This section presents our task-averaged generalization bound for Multipath MTL problem. Recall that  $\hat{\mathbf{f}} = (\hat{\mathbf{h}}, \hat{\alpha}, \hat{\phi})$  is the outcome of the ERM problem (M<sup>2</sup>TL). Observe that, if we were solving the problem with only one task, the generalization bound would depend on only one module per layer rather than the overall size of the supernet. This is because each task gets to select a single module through their pathway. In light of this, we can quantify the utilization of supernet layers as follows: Let  $\hat{K}_{\ell}$  be the number of modules utilized by the empirical solution  $\hat{\mathbf{f}}$ . Formally,  $\hat{K}_{\ell} = |\{\hat{\alpha}_t[\ell] \text{ for } t \in [T]\}|$ . The following theorem provides our guarantee in terms of Gaussian complexities of individual modules.

**Theorem 1** Suppose Assumptions 1&2 hold. Let  $\hat{\mathbf{f}}$  be the empirical solution of (M<sup>2</sup>TL). Then, with probability at least  $1 - \delta$ , the excess test risk in (2) obeys  $\mathcal{R}_{\text{M}^2\text{TL}}(\hat{\mathbf{f}})$

$$\lesssim \widetilde{\mathcal{G}}_N(\mathcal{H}) + \sum_{\ell=1}^L \sqrt{\hat{K}_{\ell}} \widetilde{\mathcal{G}}_{NT}(\Psi_{\ell}) + \sqrt{\frac{\log |\mathcal{A}|}{N} + \frac{\log(2/\delta)}{NT}}.$$

Here, the input spaces for  $\mathcal{H}$  and  $\Psi_{\ell}$  are  $\mathcal{X}_{\mathcal{H}} = \Psi_L \circ \dots \circ \Psi_1 \circ \mathcal{X}$ ,  $\mathcal{X}_{\Psi_{\ell}} = \Psi_{\ell-1} \circ \dots \circ \Psi_1 \circ \mathcal{X}$  for  $\ell > 1$ , and  $\mathcal{X}_{\Psi_1} = \mathcal{X}$ .

In Theorem 1,  $\sqrt{\frac{\log |\mathcal{A}|}{N}}$  quantifies the cost of learning the pathway and  $\widetilde{\mathcal{G}}_N(\mathcal{H})$  quantifies the cost of learning the prediction head for each task  $t \in [T]$ .  $\log |\mathcal{A}|$  dependence is standard for the discrete search space  $|\mathcal{A}|$ . The  $\widetilde{\mathcal{G}}_{NT}(\Psi_{\ell})$  terms are more interesting and reflect the benefits of MTL. The reason is that, these modules are essentially learned with

$NT$  samples rather than  $N$  samples, thus cost of representation learning is shared across tasks. The  $\sqrt{\widehat{K}_\ell}$  multiplier highlights the fact that, we only need to worry about the used modules rather than all possible  $K_\ell$  modules we could have used. In essence,  $\sum_{\ell=1}^L \sqrt{\widehat{K}_\ell} \widetilde{\mathcal{G}}_{NT}(\Psi_\ell)$  summarizes the Gaussian complexity of  $\widehat{\mathcal{G}}(\Phi_{\text{used}})$  where  $\Phi_{\text{used}}$  is the subnet of the supernet utilized by the ERM solution  $\hat{f}$ . By definition  $\widehat{\mathcal{G}}(\Phi_{\text{used}}) \leq \widehat{\mathcal{G}}(\Phi)$ . With all these in mind, Theorem 1 formalizes our earlier statement (1).

A key challenge we address in Theorem 1 is decomposing the complexity of the combined hypothesis class  $\mathcal{F}$  in (M<sup>2</sup>TL) into its building blocks  $\mathcal{A}, \mathcal{H}, (\Psi_\ell)_{\ell=1}^L$ . This is accomplished by developing Gaussian complexity chain rules inspired from the influential work of (Tripuraneni, Jordan, and Jin 2020; Maurer 2016). While this work focuses on two layer composition (prediction heads composed with a shared representation), we develop bounds to control arbitrarily long compositions of hypotheses. Accomplishing this in our multipath setting presents additional technical challenges because each task gets to choose a unique pathway. Thus, tasks don't have to contribute to the learning process of each module unlike the vanilla MTL with shared representation. Consequently, ERM solution is highly heterogeneous and some modules and tasks will be learned better than the others. Worst-case Gaussian complexity plays an important role to establish clean upper bounds in the face of this heterogeneity. In fact, in supplementary material, we provide tighter bounds in terms of empirical Gaussian complexity  $\widehat{\mathcal{G}}$ , however, they necessitate more convoluted definitions that involve the number of tasks that choose a particular module.

Finally, we note that our bound has a natural interpretation for parametric classes whose  $\log(\varepsilon$ -covering number) (i.e. metric entropy) grows with degrees of freedom as  $\text{DoF} \cdot \log(1/\varepsilon)$ . Then, Theorem 1 implies a risk bound proportional to

$\sqrt{\frac{T \cdot (\text{DoF}(\mathcal{H}) + \log |\mathcal{A}|) + \sum_{\ell=1}^L \widehat{K}_\ell \cdot \text{DoF}(\Psi_\ell)}{NT}}$ . For a neural net implementation, this means small risk as soon as total sample size  $NT$  exceeds total number of weights.

We have a few more remarks in place, discussed below.

• **Dependencies.** In Theorem 1,  $\lesssim$  suppresses dependencies on  $\log(NT)$  and  $\Gamma^L$ . The latter term arises from the exponentially growing Lipschitz constant as we compose more/deeper modules, however, it can be treated as a constant for fixed depth  $L$ . We note that such exponential depth dependence is frequent in existing generalization guarantees in deep learning literature (Golowich, Rakhlin, and Shamir 2018; Bartlett, Foster, and Telgarsky 2017; Neyshabur et al. 2018, 2017). In supplementary material, we prove that the exponential dependence can be replaced with a much stronger bound of  $\sqrt{L}$  by assuming parameterized hypothesis classes.

• **Implications for Vanilla MTL.** Observe that Vanilla MTL with single shared representation corresponds to the setting  $L = 1$  and  $K_1 = 1$ . Also supernet is simply  $\Phi = \Psi_1$  and  $\log |\mathcal{A}| = 0$ . Applying Theorem 1 to this setting with  $T$  tasks each with  $N$  samples, we obtain an excess risk upper bound of  $\widetilde{\mathcal{O}}\left(\widetilde{\mathcal{G}}_{NT}(\Phi) + \widetilde{\mathcal{G}}_N(\mathcal{H})\right)$ , where representation  $\Phi$

is trained with  $NT$  samples with input space  $\mathcal{X}$ , and task-specific heads  $h_t \in \mathcal{H}$  are trained with  $N$  samples with input space  $\Phi \circ \mathcal{X}$ . This bound recovers earlier guarantees by (Maurer, Pontil, and Romera-Paredes 2016; Tripuraneni, Jordan, and Jin 2020).

• **Unselected modules do not hurt performance.** A useful feature of our bound is its dependence on  $\Phi_{\text{used}}$  (spanned by empirical pathways) rather than full hypothesis class  $\Phi$ . This feature arises from a uniform concentration argument where we uniformly control the excess MTL risk over all potential  $\Phi_{\text{used}}$  choices. This uniform control ensures  $\widetilde{\mathcal{G}}_{NT}(\Phi_{\text{used}})$  cost for the actual solution  $\hat{f}$  and it only comes at the cost of an additional  $\sqrt{\frac{\log |\mathcal{A}|}{N}}$  term which is free (up to constant)!

• **Continuous pathways.** This work focuses on relatively simple pathways where tasks choose one module from each layer. The results can be extended to other choices of pathway sets  $\mathcal{A}$ . First, note that, as long as  $\mathcal{A}$  is a discrete set, we will naturally end up with the excess risk dependence of  $\sqrt{\frac{\log |\mathcal{A}|}{N}}$ . However, one can also consider continuous  $\alpha$ , for instance, due to relaxation of the discrete set with a simplex constraint. Such approaches are common in differentiable architecture search methods (Liu, Simonyan, and Yang 2019). In this case, each entry  $\alpha[\ell]$  can be treated as a  $K_\ell$  dimensional vector that chooses a continuous superposition of  $\ell$ 'th layer modules. Thus, the overall  $\alpha \in \mathcal{A}$  parameter would have  $\text{comp}(\mathcal{A}) = \sum_{\ell=1}^L K_\ell$  resulting in an excess risk term of  $\sqrt{\sum_{\ell=1}^L K_\ell/N}$ . Note that, these are high-level insights based on classical generalization arguments. In practice, performance can be much better than these uniform concentration based upper bounds.

• **No harm under overparameterization.** A drawback of Theorem 1 is that, it is an average-risk guarantee over  $T$  tasks. In practice, it is possible that some tasks are hurt during MTL because they are isolated or dissimilar to others (see supplementary for examples). Below, we show that, if the supernet achieves zero empirical risk, then, no task will be worse than the scenario where they are individually trained with  $N$  samples, i.e. Multipath MTL does not hurt any task.

**Lemma 1** Recall  $\hat{f}$  is the solution of (M<sup>2</sup>TL) and  $\hat{f}_t = \hat{h}_t \circ \hat{\phi}_{\hat{\alpha}_t}$  is the associated task- $t$  hypothesis. Define the excess risk of task  $t$  as  $\mathcal{R}_t(\hat{f}_t) = \mathcal{L}_t(\hat{f}_t) - \mathcal{L}_t^*$  where  $\mathcal{L}_t(f) = \mathbb{E}_{\mathcal{D}_t}[\widehat{\mathcal{L}}_t(f)]$  is the population risk of task  $t$  and  $\mathcal{L}_t^*$  is the optimal achievable test risk for task  $t$  over  $\mathcal{F}$ . With probability at least  $1 - \delta - \mathbb{P}(\widehat{\mathcal{L}}_{S_{\text{all}}}(\hat{f}) \neq 0)$ , for all tasks  $t \in [T]$ ,

$$\mathcal{R}_t(\hat{f}_t) \lesssim \widetilde{\mathcal{G}}_N(\mathcal{H}) + \sum_{\ell=1}^L \widetilde{\mathcal{G}}_N(\Psi_\ell) + \sqrt{\frac{\log(2T/\delta)}{N}}.$$

Here,  $\mathbb{P}(\widehat{\mathcal{L}}_{S_{\text{all}}}(\hat{f}) = 0)$  is the event of interpolation (zero empirical risk) under which the guarantee holds. We call this *no harm* because the bound is same as what one would get by applying union bound over  $T$  empirical risk minimizations where each task is optimized individually.

### 3.2 Transfer Learning with Optimal Pathway

Following Multipath MTL problem, in this section, we discuss guarantees for transfer learning on a supernet. Recall that  $\mathcal{A}$  is the set of pathways and our goal in (TLOP) is finding the optimal pathway  $\alpha \in \mathcal{A}$  and prediction head  $h \in \mathcal{H}_{\mathcal{T}}$  to achieve small target risk. In order to quantify the bias arising from the Multipath MTL phase, we introduce the following definition.

**Definition 2 (Supernet Bias)** Recall the definitions  $\mathcal{D}_{\mathcal{T}}$ ,  $\mathcal{H}_{\mathcal{T}}$ , and  $\mathcal{L}_{\mathcal{T}}^*$  stated in Section 2. Given a supernet  $\phi$ , we define the supernet/representation bias of  $\phi$  for a target  $\mathcal{T}$  as

$$\text{Bias}_{\mathcal{T}}(\phi) = \min_{h \in \mathcal{H}_{\mathcal{T}}, \alpha \in \mathcal{A}} \mathcal{L}_{\mathcal{T}}(h \circ \phi_{\alpha}) - \mathcal{L}_{\mathcal{T}}^*.$$

Definition 2 is a restatement of the supernet bias term in (3). Importantly, it ensures that the optimal pathway representation over  $\phi$  can not be worse than the optimal performance by  $\text{Bias}_{\mathcal{T}}(\phi)$ . Following this, we can state a generalization guarantee for transfer learning problem (TLOP).

**Theorem 2** Suppose Assumptions 1&2 hold. Let supernet  $\hat{\phi}$  be the solution of (M<sup>2</sup>TL) and  $\hat{f}_{\hat{\phi}}$  be the empirical minima of (TLOP) with respect to supernet  $\hat{\phi}$ . Then with probability at least  $1 - \delta$ ,

$$\mathcal{R}_{\text{TLOP}}(\hat{f}_{\hat{\phi}}) \lesssim \text{Bias}_{\mathcal{T}}(\hat{\phi}) + \sqrt{\frac{\log(2|\mathcal{A}|/\delta)}{M}} + \tilde{\mathcal{G}}_M(\mathcal{H}_{\mathcal{T}}),$$

where input space of  $\tilde{\mathcal{G}}_M(\mathcal{H}_{\mathcal{T}})$  is given by  $\{\hat{\phi}_{\alpha} \circ \mathcal{X} \mid \alpha \in \mathcal{A}\}$ .

Theorem 2 highlights the sample efficiency of transfer learning with optimal pathway. While the derivation is straightforward relative to Theorem 1, the key consideration is the supernet bias  $\text{Bias}_{\mathcal{T}}(\hat{\phi})$ . This term captures the excess risk in (TLOP) introduced by using  $\hat{\phi}$ . Let  $\phi^*$  be the population minima of (M<sup>2</sup>TL). Then we can define the *supernet distance* of  $\hat{\phi}$  and  $\phi^*$  by  $d_{\mathcal{T}}(\hat{\phi}; \phi^*) = \text{Bias}_{\mathcal{T}}(\hat{\phi}) - \text{Bias}_{\mathcal{T}}(\phi^*)$ . The distance measures how well the finite sample solution  $\hat{\phi}$  from (M<sup>2</sup>TL) performs compared to the optimal MTL solution  $\phi^*$ . A plausible assumption is so-called *task diversity* proposed by Chen et al. (2021); Tripuraneni, Jordan, and Jin (2020); Xu and Tewari (2021). Here, the idea (or assumption) is that, if a target task is similar to the source tasks, the distance term for target can be controlled in terms of the excess MTL risk  $\mathcal{R}_{\text{M}^2\text{TL}}(\hat{f})$  (e.g. by assuming  $d_{\mathcal{T}}(\hat{\phi}; \phi^*) \lesssim \mathcal{R}_{\text{M}^2\text{TL}}(\hat{f}) + \varepsilon$ ). Plugging in this assumption would lead to end-to-end transfer guarantees by integrating Theorems 1 and 2, and we extend the formal analysis to appendix. However, as discussed in Theorem 4, in multipath setting, the problem is a lot more intricate because source tasks can choose totally different task-specific representations making such assumptions unrealistic. In contrast, Theorem 4 establishes concrete guarantees by probabilistically relating target and source distributions. Finally,  $\text{Bias}_{\mathcal{T}}(\phi^*)$  term is unavoidable, however, similar to  $d_{\mathcal{T}}(\hat{\phi}; \phi^*)$ , it will be small as long as source and target tasks benefit from a shared supernet at the population level.

### 4 Guarantees for Linear Representations

As a concrete instantiation of Multipath MTL, consider a linear representation learning problem where each module  $\psi_{\ell}^k$  applies matrix multiplications parameterized by  $\mathbf{B}_{\ell}^k$  with dimensions  $p_{\ell} \times p_{\ell-1}$ :  $\psi_{\ell}^k(\mathbf{x}) = \mathbf{B}_{\ell}^k \mathbf{x}$ . Here  $p_{\ell}$  are module dimensions with input dimension  $p_0 = p$  and output dimension  $p_L$ . Given a path  $\alpha$ , we obtain the linear representation  $\mathbf{B}_{\alpha} = \prod_{\ell=1}^L \mathbf{B}_{\ell}^{\alpha[\ell]} \in \mathbb{R}^{p_L \times p}$  where  $p_L$  is the number of rows of the final module  $\mathbf{B}_L^{\alpha[L]}$ . When  $p_L \ll p$ ,  $\mathbf{B}_{\alpha}$  is a fat matrix that projects  $\mathbf{x} \in \mathbb{R}^p$  onto a lower dimensional subspace. This way, during few-shot adaptation, we only need to train  $p_L \ll p$  parameters with features  $\mathbf{B}_{\alpha} \mathbf{x}$ . This is also the central idea in several works on linear meta-learning (Kong et al. 2020a; Sun et al. 2021; Bouniot et al. 2020; Tripuraneni, Jin, and Jordan 2021) which focus on a single linear representation. Our discussion within this section extends these results to the Multipath MTL setting.

Denote  $\mathbf{f} = \{((\mathbf{B}_{\ell}^k)_{k=1}^{K_{\ell}})_{\ell=1}^L, (\mathbf{h}_t, \alpha_t)_{t=1}^T\}$  where  $\mathbf{h}_t \in \mathbb{R}^{p_L}$  are linear prediction heads. Let  $\mathcal{F}$  be the search space associated with  $\mathbf{f}$ . Follow the similar setting as in Section 2 and let  $\mathcal{X} \subset \mathbb{R}^p$ . Given dataset  $\mathcal{S}_{\text{all}} = (\mathcal{S}_t)_{t=1}^T$ , we study

$$\hat{\mathbf{f}} = \min_{\mathbf{f} \in \mathcal{F}} \hat{\mathcal{L}}_{\mathcal{S}_{\text{all}}}(\mathbf{f}) := \frac{1}{NT} \sum_{t=1}^T \sum_{i=1}^N (y_{ti} - \mathbf{h}_t^{\top} \mathbf{B}_{\alpha_t} \mathbf{x}_{ti})^2. \quad (4)$$

Let  $\mathcal{B}^p(r) \subset \mathbb{R}^p$  be the Euclidean ball of radius  $r$ . To proceed, we make the following assumption for a constant  $C \geq 1$ .

**Assumption 3** For all  $\ell \in [L]$ ,  $\Psi_{\ell}$  is the set of matrices with operator norm bounded by  $C$  and  $\mathcal{H} = \mathcal{B}^{p_L}(C)$ .

The result below is a variation of Theorem 1 where the bound is refined for linear representations (with finite parameters).

**Theorem 3** Suppose Assumptions 2&3 hold, and input set  $\mathcal{X} \subset \mathcal{B}^p(R)$  for a constant  $R > 0$ . Then, with probability at least  $1 - \delta$ ,

$$\mathcal{R}_{\text{M}^2\text{TL}}(\hat{\mathbf{f}}) \lesssim \sqrt{\frac{L \cdot \text{DoF}(\mathcal{F})}{NT}} + \sqrt{\frac{\log |\mathcal{A}|}{N} + \frac{\log(2/\delta)}{NT}},$$

where  $\text{DoF}(\mathcal{F}) = T \cdot p_L + \sum_{\ell=1}^L K_{\ell} \cdot p_{\ell} \cdot p_{\ell-1}$  is the total number of trainable parameters in  $\mathcal{F}$ .

We note that Theorem 3 can be stated more generally for neural nets by placing ReLU activations between layers. Here  $\lesssim$  subsumes the logarithmic dependencies, and the sample complexity has linear dependence on  $L$  (rather than exponential dependence as in Thm 1). In essence, it implies small task-averaged excess risk as soon as total sample size  $\gtrsim$  total number of weights.

While flexible, this result does not guarantee that  $\hat{\mathbf{f}}$  can benefit transfer learning for a new task. To proceed, we introduce additional assumptions under which we can guarantee the success of (TLOP). The first assumption is a realizability condition that guarantees tasks share same supernet representation (so that supernet bias is small).

**Assumption 4 (A)** Task datasets are generated from a planted model  $(\mathbf{x}_t, y_t) \sim \mathcal{D}_t$  where  $y_t = \mathbf{x}_t^{\top} \boldsymbol{\theta}_t^* + z_t$  where

$\mathbf{x}_t, z_t$  are zero mean,  $\mathcal{O}(1)$ -subgaussian and  $\mathbb{E}[\mathbf{x}_t \mathbf{x}_t^\top] = \mathbf{I}_p$ .  
**(B) Task vectors are generated according to ground-truth supernet**  $\mathbf{f}^* = \{((\bar{\mathbf{B}}_\ell^k)_{k=1}^{K_\ell})_{\ell=1}^L, (\bar{\mathbf{h}}_t, \bar{\alpha}_t)_{t=1}^T\}$  so that  $\boldsymbol{\theta}_t^* = \bar{\mathbf{B}}_{\bar{\alpha}_t}^\top \bar{\mathbf{h}}_t$ .  $\mathbf{f}^*$  is normalized so that  $\|\bar{\mathbf{B}}_\ell^k\| = \|\bar{\mathbf{h}}_t\| = 1$ .

Our second assumption is a task diversity condition adapted from (Tripuraneni, Jin, and Jordan 2021; Kong et al. 2020b) that facilitates the identifiability of the ground truth supernet.

**Assumption 5 (Diversity during MTL)** Cluster the tasks by their pathways via  $\mathbf{H}_\alpha = \{\mathbf{h}_t \mid \bar{\alpha}_t = \alpha\}$ . Define cluster population  $\gamma_\alpha = |\mathbf{H}_\alpha|/p_L$  and covariance  $\boldsymbol{\Sigma}_\alpha = \gamma_\alpha^{-1} \sum_{\mathbf{h} \in \mathbf{H}_\alpha} \mathbf{h} \mathbf{h}^\top$ . For a proper constant  $c > 0$  and for all pathways  $\alpha$  we have  $\boldsymbol{\Sigma}_\alpha \succeq c \mathbf{I}_{p_L}$ .

Verbally, this condition requires that, if a pathway is chosen by a source task, that pathway should contain diverse tasks so that (M<sup>2</sup>TL) phase can learn a good representation that can benefit transfer learning. However, this definition is flexible in the sense that pathways can still have sophisticated interactions/intersections and we don't assume anything for the pathways that are not chosen by source. We also have the challenge that, some pathways can be a lot more populated than others and target task might suffer from poor MTL representation quality over less populated pathways. The following assumption is key to overcoming this issue by enforcing a distributional prior on the target task pathway so that its pathway is similar to the source tasks in average.

**Assumption 6 (Distribution of target task)** Draw  $\alpha_T$  uniformly at random from source pathways  $(\bar{\alpha}_t)_{t=1}^T$ . Target task is distributed as in Assumption 4(A) with pathway  $\alpha_T$  and  $\boldsymbol{\theta}_T^* = \bar{\mathbf{B}}_{\alpha_T}^\top \mathbf{h}_T$  with  $\|\mathbf{h}_T\| = 1$ .

With these assumptions, we have the following result that guarantees end-to-end multipath learning ((M<sup>2</sup>TL) phase followed by (TLOP) using MTL representation).

**Theorem 4** Suppose Assumptions 3–6 hold and  $\ell(\hat{y}, y) = (y - \hat{y})^2$ . Additionally assume input set  $\mathcal{X} \subset \mathcal{B}^p(R)$  for a constant  $R > 0$  and  $\mathcal{H}_T \subset \mathbb{R}^{p_L}$ . Solve MTL problem (M<sup>2</sup>TL) with the knowledge of ground-truth pathways  $(\bar{\alpha}_t)_{t=1}^T$  to obtain a supernet  $\hat{\phi}$  and  $NT \gtrsim \text{DoF}(\mathcal{F}) \log(NT)$ . Solve transfer learning problem (TLOP) with  $\hat{\phi}$  to obtain a target hypothesis  $\hat{f}_\phi$ . Then, with probability at least  $1 - 3e^{-cM} - \delta$ , path-averaged excess target risk (3) obeys  $\mathbb{E}_{\alpha_T}[\mathcal{R}_{\text{TLOP}}(\hat{f}_\phi)]$

$$\lesssim p_L \sqrt{\frac{L \cdot \text{DoF}(\mathcal{F}) + \log(8/\delta)}{NT}} + \frac{p_L}{M} + \sqrt{\frac{\log(8|\mathcal{A}|/\delta)}{M}}.$$

Here  $\text{DoF}(\mathcal{F}) = T \cdot p_L + \sum_{\ell=1}^L K_\ell \cdot p_\ell \cdot p_{\ell-1}$ , and  $\mathbb{E}_{\alpha_T}$  denotes the expectation over the random target pathways.

In words, this result controls the target risk in terms of the sample size of the target task and sample size during multi-task representation learning, and provides a concrete instantiation of discussion following Theorem 2. In Theorem 9 in appendix, we provide a tighter bound for expected transfer risk when linear head  $\mathbf{h}_T$  is uniformly drawn from the unit sphere. The primary challenge in our work compared to related vanilla MTL results by (Tripuraneni, Jin, and Jordan 2021; Du et al. 2020; Kong et al. 2020b) is the fact that,

we deal with exponentially many pathway representations many of which may be low quality. Assumption 6 allows us to convert task-averaged MTL risk into a transfer learning guarantee over a random pathway. Finally, Theorem 4 assumes that source pathways are known during MTL phase. In Appendix E, we show that this assumption is indeed necessary: Otherwise, one can construct scenarios where (M<sup>2</sup>TL) problem admits an alternative solution  $\hat{f}$  with optimal MTL risk but the resulting supernet  $\hat{\phi}$  achieves poor target risk. Supplementary material discusses this challenge and identifies additional conditions that make ground-truth pathways uniquely identifiable when we solve (M<sup>2</sup>TL).

## 5 Insights from Hierarchical Representations

We now discuss the special two-layer supernet structure depicted in Figure 1b. This setting groups tasks into  $K := K_2$  clusters and first layer module is shared across all tasks ( $K_1 = 1$ ). Ignoring first layer, pathway  $\alpha_t \in [K]$  becomes the clustering assignment for task  $t$ . Applying Theorem 1, we obtain a generalization bound of

$$\mathcal{R}_{\text{M}^2\text{TL}}(\hat{f}) \lesssim \tilde{\mathcal{G}}_{NT}(\Psi_1) + \sqrt{K} \tilde{\mathcal{G}}_{NT}(\Psi_2) + \tilde{\mathcal{G}}_N(\mathcal{H}) + \sqrt{\frac{\log K}{N}}.$$

Here,  $\psi_1 \in \Psi_1$  is the shared first layer module,  $\psi_2^k \in \Psi_2$  is the module assigned to cluster  $k \in [K]$  that personalizes its representation, and we have  $|\mathcal{A}| = K$ . To provide further insights, let us focus on linear representations with the notation of Section 4:  $\psi_1(\mathbf{x}) = \mathbf{B}_1 \mathbf{x}$ ,  $\psi_2^k(\mathbf{x}') = \mathbf{B}_2^k \mathbf{x}'$ , and  $h_t(\mathbf{x}'') = \mathbf{h}_t^\top \mathbf{x}''$  with dimensions  $\mathbf{B}_1 \in \mathbb{R}^{R \times p}$ ,  $\mathbf{B}_2^k \in \mathbb{R}^{r \times R}$ ,  $\mathbf{h}_t \in \mathbb{R}^r$  and  $r \leq R \leq p$ . Our bound now takes the form

$$\mathcal{R}_{\text{M}^2\text{TL}}(\hat{f}) \lesssim \sqrt{\frac{Rp + KrR + T(r + \log K)}{NT}},$$

where  $Rp$  and  $KrR$  are the number of parameters in supernet layers 1 and 2, and  $(r + \log K)/N$  is the cost of learning pathway and prediction head per task. Let us contrast this to the shallow MTL approaches with 1-layer supernet.

- **Vanilla MTL:** Learn  $\mathbf{B}_1 \in \mathbb{R}^{R \times p}$  and learn larger prediction heads  $\mathbf{h}_t^V \in \mathbb{R}^R$  (no clustering needed).
- **Cluster MTL:** Learn larger cluster modules  $\mathbf{B}_2^{C,k} \in \mathbb{R}^{r \times p}$ , and learn pathway  $\alpha_t$  and head  $\mathbf{h}_t \in \mathbb{R}^r$  (no  $\mathbf{B}_1$  needed).

**Experimental Insights.** Before providing a theoretical comparison, let us discuss the experimental results where we compare these three approaches in a realizable dataset generated according to Figure 1b. Specifically, we generate  $\bar{\mathbf{B}}_1$  and  $\{\bar{\mathbf{B}}_2^k\}_{k=1}^K$  with orthonormal rows uniformly at random independently. We also generate  $\bar{\mathbf{h}}_t$  uniformly at random over the unit sphere independently. Let  $\bar{\alpha}_t$  be the cluster assignment of task  $t$  where each cluster has same size/number of tasks with  $T = T/K$  tasks. The distribution  $\mathcal{D}_t$  associated with task  $t$  is generated as

$$y = \mathbf{x}^\top \boldsymbol{\theta}_t^* \quad \text{where} \quad \boldsymbol{\theta}_t^* = (\bar{\mathbf{h}}_t^\top \bar{\mathbf{B}}_2^{\alpha_t} \bar{\mathbf{B}}_1)^\top, \quad \mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_p),$$

without label noise. We evaluate and present results from two scenarios where cluster assignment of each task  $\bar{\alpha}_t$  is known (Figure 2) or not (Figure 3). MTL, Cluster-MTL and

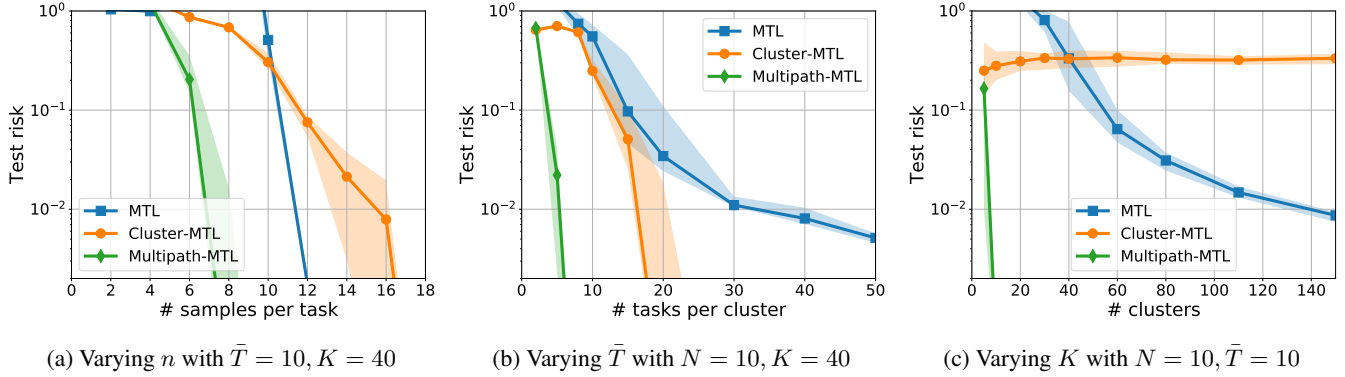


Figure 2: We compare the sample complexity of MTL, Cluster-MTL and Multipath-MTL in a noiseless linear regression setting. For each figure, we fix two of the configurations and vary the other one. We find that Multipath-MTL is superior to both baselines as predicted by our theory. The solid curves are the median risk and the shaded regions highlight the first and third quantile risks. Each marker is obtained by averaging 20 independent realizations.

Multipath-MTL labels corresponds to our single representation, clustering and hierarchical MTL strategies respectively, in the figures.

In Figure 2, we solve MTL problems with the knowledge of clustering  $\bar{\alpha}_t$ . We set ambient dimension  $p = 32$ , shared embedding  $R = 8$ , and cluster embeddings  $r = 2$ . We consider a base configuration of  $K = 40$  clusters,  $\bar{T} = T/K = 10$  tasks per cluster and  $N = 10$  samples per task (see supplementary material for further details). Figure 2 compares the performance of three approaches for the task-averaged MTL test risk and demonstrates consistent benefits of Multipath MTL for varying  $K, \bar{T}, N$ .

We also consider the setting where  $\bar{\alpha}_t, t \in [T]$  are unknown during training. Set  $p = 128, R = 32$  and  $r = 2$ , and fix number of clusters  $K = 50$  and cluster size  $\bar{T} = 10$ . In this experiment, instead of using the ground truth clustering  $\bar{\alpha}_t$ , we also learn the clustering assignment  $\hat{\alpha}_t$  for each task. As we discussed and visualized in supplementary material, it is not easy to cluster random tasks even with the hindsight knowledge of task vectors  $\theta_t^*$ . To overcome this issue, we add correlation between tasks in the same cluster. Specifically, generate the prediction head by  $\bar{h}_t = \gamma \bar{h}^k + (1 - \gamma) \bar{h}_t$  where  $\bar{h}^k, \bar{h}_t$  are random unit vectors corresponding to the cluster  $k$  and task  $t$  (assuming  $\bar{\alpha}_t = k$ ). To cluster tasks, we first run vanilla MTL and learn the shared representation  $\hat{B}_1$  and heads  $(\hat{h}_t^V)_{t=1}^T$ . Next build task vector estimates by  $\hat{\theta}_t := \hat{B}_1^\top \hat{h}_t^V$ , and get  $T \times T$  task similarity matrix using Euclidean distance metric. Applying standard  $K$ -means clustering to it provides a clustering assignment  $\hat{\alpha}_t$ . In the experiment, we set  $\gamma = 0.6$  to make sure hindsight knowledge of  $\theta_t^*$  is sufficient to correctly cluster all tasks. Results are presented in Figure 3, where solid curves are solving MTL with ground truth  $\bar{\alpha}_t$  while dashed curves are using  $\hat{\alpha}_t$ . We observe that when given enough samples ( $N \geq 60$ ), all tasks are grouped correctly even if the MTL risk is not zero. More importantly, Multipath MTL does outperform both vanilla MTL and cluster MTL even when the clustering is not fully correct.

**Understanding the benefits of Multipath MTL.** Naturally, superior numerical performance of Multipath MTL in Figure 2&3 partly stems from the hierarchical dataset model we

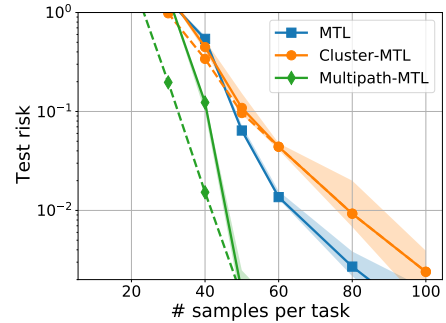


Figure 3: We group the  $T = 500$  tasks into  $K = 50$  clusters and compare the sample complexity of different MTL strategies. Given different sample size, we cluster tasks based on the trained MTL model and solve Cluster-/Multipath-MTL based on the assigned clusters. Solid curves are results using ground truth cluster knowledge  $\bar{\alpha}_t$  and dashed are using the learned clustering  $\hat{\alpha}_t$ . Experimental setting follows the same setting as in Figure 2.

study. This model will also shed light on shortcomings of 1-layer supernet drawing from our theoretical predictions. First, observe that all three baselines are exactly specified: We use the smallest model sizes that capture the ground-truth model so that they can achieve zero test risk as  $N, K, T$  grows. For instance, Vanilla MTL achieves zero risk by setting  $B_1 = \bar{B}_1, h_t = (\bar{B}_2^{\alpha_t})^\top \bar{h}_t$  and cluster MTL achieves zero risk by setting  $B_2^{C,k} = \bar{B}_2^k \bar{B}_1, h_t = \bar{h}_t$ . Thus, the benefit of Multipath MTL arises from stronger weight sharing across tasks that reduces test risk. In light of Sec. 4, the generalization risks of these approaches can be bounded as  $\sqrt{\text{DoF}(\mathcal{F})/NT}$  where Number-of-Parameters compare as **Vanilla:**  $Rp + TR$ , **Cluster:**  $Krp + Tr$ , **Multipath:**  $Rp + KrR + Tr$ . From this, it can be seen that Multipath is never worse than the others as long as  $Kr \geq R$  and  $\bar{T} = T/K \geq r$ . These conditions hold under the assumption that multipath model is of minimal size: Otherwise, there would be a strictly smaller zero-risk model by setting  $R \leftarrow Kr$  and  $r \leftarrow \bar{T}$ .

Conversely, Multipath shines in the regime  $Kr \gg R$  or  $\bar{T} \gg r$ . As  $\frac{Kr}{R}, \frac{p}{R} \rightarrow \infty$ , Multipath strictly outperforms Cluster MTL. This arises from a *cluster diversity* phe-

nomenon that connects to the *task diversity* notions of prior art. In essence, since  $r$ -dimensional clusters lie on a shared  $R$  dimensional space, as we add more clusters beyond  $Kr \geq R$ , they will collaboratively estimate the shared subspace which in turn helps estimating their local subspaces by projecting them onto the shared one. As  $\frac{T}{r}, \frac{R}{r} \rightarrow \infty$ , Multipath strictly outperforms Vanilla MTL.  $\frac{T}{r}$  is needed to ensure that there is enough task diversity within each cluster to estimate its local subspace. Finally,  $\frac{R}{r}$  ratio is the few-shot learning benefit of clustering over Vanilla MTL. The prediction heads of vanilla MTL is larger which necessitates a larger  $N$ , at the minimum  $N \geq R$ . Whereas Multipath works with as little as  $N \geq r$ . The same argument also implies that clustering/hierarchy would also enable better transfer learning.

## 6 Related Work

Our work is related to a large body of literature spanning efficient architectures and statistical guarantees for MTL, representation learning, task similarity, and subspace clustering.

• **Multitask Representation Learning.** While MTL problems admit multiple approaches, an important idea is building shared representations to embed tasks in a low-dimensional space (Zhang and Yang 2021; Thrun and Pratt 2012; Wang, Kolar, and Srebro 2016; Baxter 2000). After identifying this low-dimensional representation, new tasks can be learned in a sample efficient fashion inline with the benefits of deep representations in modern ML applications. While most earlier works focus on linear models, (Maurer, Pontil, and Romera-Paredes 2016) provides guarantees for general hypothesis classes through empirical process theory improving over (Baxter 2000). More recently, there is a growing line of work on multitask representations that spans tighter sample complexity analysis (Garg and Liang 2020; Hanneke and Kpotufe 2020; Du et al. 2020; Kong et al. 2020b; Xu and Tewari 2021; Lu, Huang, and Du 2021), convergence guarantees (Balcan, Khodak, and Talwalkar 2019; Khodak, Balcan, and Talwalkar 2019; Collins et al. 2022; Ji et al. 2020; Collins et al. 2021; Wu, Zhang, and Ré 2020), lifelong learning (Xu and Tewari 2022; Li et al. 2022), and decision making problems (Yang et al. 2020; Qin et al. 2022; Cheng et al. 2022; Sodhani, Zhang, and Pineau 2021). Closest to our work is (Tripuraneni, Jin, and Jordan 2021) which provides tighter sample complexity guarantees compared to (Maurer, Pontil, and Romera-Paredes 2016). Our problem formulation generalizes prior work (that is mostly limited to single shared representation) by allowing deep compositional representations computed along supernet pathways. To overcome the associated technical challenges, we develop multilayer chain rules for Gaussian Complexity, introduce new notions to assess the quality of supernet representations, and develop new theory for linear representations.

• **Quantifying Task Similarity and Clustering.** We note that task similarity and clustering has been studied by (Shui et al. 2019; Nguyen, Do, and Carneiro 2021; Zhou et al. 2020; Fifty et al. 2021; Kumar and Daume III 2012; Kang, Grauman, and Sha 2011; Aribandi et al. 2021; Zamir et al. 2018) however these works do not come with comparable statistical guarantees. Leveraging relations between tasks are

explored even more broadly (Zhuang et al. 2020; Achille et al. 2021). Our experiments on linear Multipath MTL connects well with the broader subspace clustering literature (Vidal 2011; Parsons, Haque, and Liu 2004; Elhamifar and Vidal 2013). Specifically, each learning task  $\theta_i$  can be viewed as a point on a high-dimensional subspace. Multipath MTL aims to cluster these points into smaller subspaces that correspond to task-specific representations. Our challenge is that we only get to see the points through the associated datasets.

• **ML Architectures and Systems.** While traditional ML models tend to be good at a handful of tasks, next-generation of neural architectures are expected to excel at a diverse range of tasks while allowing for multiple input modalities. To this aim, task-specific representations can help address both computational and data efficiency challenges. Recent works (Ramesh and Chaudhari 2021a; Shu et al. 2021; Ramesh and Chaudhari 2021b; Fifty et al. 2021; Yao et al. 2019; Vuorio et al. 2019; Mansour et al. 2020; Tan et al. 2022; Ghosh et al. 2020; Collins et al. 2021) propose hierarchical/clustering approaches to group tasks in terms of their similarities, (Qin et al. 2020; Ye, Zha, and Ren 2022; Gupta et al. 2022; Asai et al. 2022; He et al. 2022) focus on training mixture-of-experts (MoE) models, and similar to the pathways (Strezoski, Noord, and Worring 2019; Rosenbaum, Klinger, and Riemer 2017; Chen, Gu, and Fu 2021; Ma et al. 2019) study on task routing. In the context of lifelong learning, PathNet, PackNet (Fernando et al. 2017; Mallya and Lazebnik 2018) and many other existing methods (Parisi et al. 2019; Mallya, Davis, and Lazebnik 2018; Hung et al. 2019; Wortsman et al. 2020; Cheung et al. 2019) propose to embed many tasks into the same network to facilitate sample/compute efficiency. PathNet as well as SNR (Ma et al. 2019) propose methods to identify pathways/routes for individual tasks and efficiently compute them over the conditional subnetwork. With the advent of large language models, conditional computation paradigm is witnessing a growing interest with architectural innovations such as muNet, GShard, Pathways, and PaLM (Gesmundo and Dean 2022a,b; Barham et al. 2022; Dean 2021; Lepikhin et al. 2020; Chowdhery et al. 2022; Driess et al. 2023) and provide a strong motivation for theoretically-grounded Multipath MTL methods.

## 7 Discussion

This work explored novel multitask learning problems which allow for task-specific representations that are computed along pathways of a large supernet. We established generalization bounds under a general setting which proved insightful when specialized to linear or hierarchical representations. We believe there are multiple exciting directions to explore. First, it is desirable to develop a stronger control over the generalization risk of specific groups of tasks. Our Lemma 1 is a step in this direction. Second, what are risk upper/lower bounds for Multipath MTL as we vary the depth and width of the supernet graph? Discussion in Section 5 falls under this question where we demonstrate the sample complexity benefits of Multipath MTL over traditional MTL approaches. Finally, following experiments in Section 5, can we establish similar provable guarantees for computationally-efficient algorithms (e.g. method of moments, gradient descent)?



## Acknowledgements

Authors would like to thank Zhe Zhao for helpful discussions and pointing out related works. This work was supported in part by the NSF grants CCF-2046816 and CCF-2212426, Google Research Scholar award, and Army Research Office grant W911NF2110312.

## References

- Achille, A.; Paolini, G.; Mbeng, G.; and Soatto, S. 2021. The information complexity of learning tasks, their structure and their distance. *Information and Inference: A Journal of the IMA*, 10(1): 51–72.
- Aribandi, V.; Tay, Y.; Schuster, T.; Rao, J.; Zheng, H. S.; Mehta, S. V.; Zhuang, H.; Tran, V. Q.; Bahri, D.; Ni, J.; et al. 2021. Ext5: Towards extreme multi-task scaling for transfer learning. *arXiv preprint arXiv:2111.10952*.
- Asai, A.; Salehi, M.; Peters, M. E.; and Hajjishirzi, H. 2022. Attentional Mixtures of Soft Prompt Tuning for Parameter-efficient Multi-task Knowledge Sharing. *arXiv preprint arXiv:2205.11961*.
- Balcan, M.-F.; Khodak, M.; and Talwalkar, A. 2019. Provable guarantees for gradient-based meta-learning. In *International Conference on Machine Learning*, 424–433. PMLR.
- Barham, P.; Chowdhery, A.; Dean, J.; Ghemawat, S.; Hand, S.; Hurt, D.; Isard, M.; Lim, H.; Pang, R.; Roy, S.; et al. 2022. Pathways: Asynchronous distributed dataflow for ML. *Proceedings of Machine Learning and Systems*, 4: 430–449.
- Bartlett, P. L.; Foster, D. J.; and Telgarsky, M. J. 2017. Spectrally-normalized margin bounds for neural networks. In *Advances in Neural Information Processing Systems*, 6241–6250.
- Baxter, J. 2000. A model of inductive bias learning. *Journal of artificial intelligence research*, 12: 149–198.
- Bouniot, Q.; Redko, I.; Audigier, R.; Loesch, A.; Zotkin, Y.; and Habrard, A. 2020. Towards better understanding meta-learning methods through multi-task representation learning theory. *arXiv preprint arXiv:2010.01992*.
- Brown, T.; Mann, B.; Ryder, N.; Subbiah, M.; Kaplan, J. D.; Dhariwal, P.; Neelakantan, A.; Shyam, P.; Sastry, G.; Askell, A.; et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33: 1877–1901.
- Caruana, R. 1997. Multitask learning. *Machine learning*, 28(1): 41–75.
- Chen, S.; Crammer, K.; He, H.; Roth, D.; and Su, W. J. 2021. Weighted Training for Cross-Task Learning. *arXiv preprint arXiv:2105.14095*.
- Chen, X.; Gu, X.; and Fu, L. 2021. Boosting share routing for multi-task learning. In *Companion Proceedings of the Web Conference 2021*, 372–379.
- Cheng, Y.; Feng, S.; Yang, J.; Zhang, H.; and Liang, Y. 2022. Provable benefit of multitask representation learning in reinforcement learning. *arXiv preprint arXiv:2206.05900*.
- Cheung, B.; Terekhov, A.; Chen, Y.; Agrawal, P.; and Olshausen, B. 2019. Superposition of many models into one. *Advances in neural information processing systems*, 32.
- Chowdhery, A.; Narang, S.; Devlin, J.; Bosma, M.; Mishra, G.; Roberts, A.; Barham, P.; Chung, H. W.; Sutton, C.; Gehrmann, S.; et al. 2022. Palm: Scaling language modeling with pathways. *arXiv preprint arXiv:2204.02311*.
- Collins, L.; Hassani, H.; Mokhtari, A.; and Shakkottai, S. 2021. Exploiting shared representations for personalized federated learning. In *International Conference on Machine Learning*, 2089–2099. PMLR.
- Collins, L.; Mokhtari, A.; Oh, S.; and Shakkottai, S. 2022. MAML and ANIL provably learn representations. *arXiv preprint arXiv:2202.03483*.
- Dean, J. 2021. Introducing Pathways: A next-generation AI architecture. <https://blog.google/technology/ai/introducing-pathways-next-generation-ai-architecture/>, *Google AI Blog*.
- Deng, J.; Dong, W.; Socher, R.; Li, L.-J.; Li, K.; and Fei-Fei, L. 2009. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, 248–255. Ieee.
- Driess, D.; Xia, F.; Sajjadi, M. S. M.; Lynch, C.; Chowdhery, A.; Ichter, B.; Wahid, A.; Tompson, J.; Vuong, Q.; Yu, T.; Huang, W.; Chebotar, Y.; Sermanet, P.; Duckworth, D.; Levine, S.; Vanhoucke, V.; Hausman, K.; Toussaint, M.; Greff, K.; Zeng, A.; Mordatch, I.; and Florence, P. 2023. PaLM-E: An Embodied Multimodal Language Model. In *arXiv preprint arXiv:2303.03378*.
- Du, S. S.; Hu, W.; Kakade, S. M.; Lee, J. D.; and Lei, Q. 2020. Few-shot learning via learning the representation, provably. *arXiv preprint arXiv:2002.09434*.
- Elhamifar, E.; and Vidal, R. 2013. Sparse subspace clustering: Algorithm, theory, and applications. *IEEE transactions on pattern analysis and machine intelligence*, 35(11): 2765–2781.
- Fernando, C.; Banarse, D.; Blundell, C.; Zwols, Y.; Ha, D.; Rusu, A. A.; Pritzel, A.; and Wierstra, D. 2017. Pathnet: Evolution channels gradient descent in super neural networks. *arXiv preprint arXiv:1701.08734*.
- Fifty, C.; Amid, E.; Zhao, Z.; Yu, T.; Anil, R.; and Finn, C. 2021. Efficiently identifying task groupings for multi-task learning. *Advances in Neural Information Processing Systems*, 34: 27503–27516.
- Garg, S.; and Liang, Y. 2020. Functional regularization for representation learning: A unified theoretical perspective. *Advances in Neural Information Processing Systems*, 33: 17187–17199.
- Gesmundo, A.; and Dean, J. 2022a. An Evolutionary Approach to Dynamic Introduction of Tasks in Large-scale Multitask Learning Systems. *arXiv preprint arXiv:2205.12755*.
- Gesmundo, A.; and Dean, J. 2022b. muNet: Evolving Pre-trained Deep Neural Networks into Scalable Auto-tuning Multitask Systems. *arXiv preprint arXiv:2205.10937*.
- Ghosh, A.; Chung, J.; Yin, D.; and Ramchandran, K. 2020. An efficient framework for clustered federated learning. *Advances in Neural Information Processing Systems*, 33: 19586–19597.
- Golowich, N.; Rakhlin, A.; and Shamir, O. 2018. Size-independent sample complexity of neural networks. In *Conference On Learning Theory*, 297–299. PMLR.

- Gupta, S.; Mukherjee, S.; Subudhi, K.; Gonzalez, E.; Jose, D.; Awadallah, A. H.; and Gao, J. 2022. Sparsely activated mixture-of-experts are robust multi-task learners. *arXiv preprint arXiv:2204.07689*.
- Hanneke, S.; and Kpotufe, S. 2020. A no-free-lunch theorem for multitask learning. *arXiv preprint arXiv:2006.15785*.
- He, C.; Zheng, S.; Zhang, A.; Karypis, G.; Chilimbi, T.; Soltanolkotabi, M.; and Avestimehr, S. 2022. SMILE: Scaling Mixture-of-Experts with Efficient Bi-level Routing. *arXiv preprint arXiv:2212.05191*.
- Hung, C.-Y.; Tu, C.-H.; Wu, C.-E.; Chen, C.-H.; Chan, Y.-M.; and Chen, C.-S. 2019. Compacting, picking and growing for unforgetting continual learning. *Advances in Neural Information Processing Systems*, 32.
- Ji, K.; Lee, J. D.; Liang, Y.; and Poor, H. V. 2020. Convergence of meta-learning with task-specific adaptation over partial parameters. *Advances in Neural Information Processing Systems*, 33: 11490–11500.
- Ji, Z.; and Telgarsky, M. 2018. Gradient descent aligns the layers of deep linear networks. *arXiv preprint arXiv:1810.02032*.
- Kang, Z.; Grauman, K.; and Sha, F. 2011. Learning with whom to share in multi-task feature learning. In *ICML*.
- Khodak, M.; Balcan, M.-F. F.; and Talwalkar, A. S. 2019. Adaptive gradient-based meta-learning methods. *Advances in Neural Information Processing Systems*, 32.
- Kong, W.; Somani, R.; Kakade, S.; and Oh, S. 2020a. Robust meta-learning for mixed linear regression with small batches. *Advances in neural information processing systems*, 33: 4683–4696.
- Kong, W.; Somani, R.; Song, Z.; Kakade, S.; and Oh, S. 2020b. Meta-learning for mixed linear regression. In *International Conference on Machine Learning*, 5394–5404. PMLR.
- Kumar, A.; and Daume III, H. 2012. Learning task grouping and overlap in multi-task learning. *arXiv preprint arXiv:1206.6417*.
- Lepikhin, D.; Lee, H.; Xu, Y.; Chen, D.; Firat, O.; Huang, Y.; Krikun, M.; Shazeer, N.; and Chen, Z. 2020. Gshard: Scaling giant models with conditional computation and automatic sharding. *arXiv preprint arXiv:2006.16668*.
- Li, Y.; Li, M.; Asif, M. S.; and Oymak, S. 2022. Provable and Efficient Continual Representation Learning. *arXiv preprint arXiv:2203.02026*.
- Liu, H.; Simonyan, K.; and Yang, Y. 2019. Darts: Differentiable architecture search. *ICLR*.
- Lu, R.; Huang, G.; and Du, S. S. 2021. On the power of multitask representation learning in linear mdp. *arXiv preprint arXiv:2106.08053*.
- Ma, J.; Zhao, Z.; Chen, J.; Li, A.; Hong, L.; and Chi, E. H. 2019. Snr: Sub-network routing for flexible parameter sharing in multi-task learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, 216–223.
- Mallya, A.; Davis, D.; and Lazebnik, S. 2018. Piggyback: Adapting a single network to multiple tasks by learning to mask weights. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 67–82.
- Mallya, A.; and Lazebnik, S. 2018. Packnet: Adding multiple tasks to a single network by iterative pruning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 7765–7773.
- Mansour, Y.; Mohri, M.; Ro, J.; and Suresh, A. T. 2020. Three approaches for personalization with applications to federated learning. *arXiv preprint arXiv:2002.10619*.
- Maurer, A. 2016. A chain rule for the expected suprema of Gaussian processes. *Theoretical Computer Science*, 650: 109–122.
- Maurer, A.; Pontil, M.; and Romera-Paredes, B. 2016. The benefit of multitask representation learning. *Journal of Machine Learning Research*, 17(81): 1–32.
- Mohri, M.; Rostamizadeh, A.; and Talwalkar, A. 2018. *Foundations of machine learning*. MIT press.
- Neyshabur, B.; Bhojanapalli, S.; McAllester, D.; and Srebro, N. 2017. Exploring generalization in deep learning. *Advances in neural information processing systems*, 30.
- Neyshabur, B.; Li, Z.; Bhojanapalli, S.; LeCun, Y.; and Srebro, N. 2018. Towards understanding the role of over-parametrization in generalization of neural networks. *arXiv preprint arXiv:1805.12076*.
- Nguyen, C.; Do, T.-T.; and Carneiro, G. 2021. Similarity of classification tasks. *arXiv preprint arXiv:2101.11201*.
- Oymak, S. 2018. Learning compact neural networks with regularization. In *International Conference on Machine Learning*, 3966–3975. PMLR.
- Parisi, G. I.; Kemker, R.; Part, J. L.; Kanan, C.; and Wermter, S. 2019. Continual lifelong learning with neural networks: A review. *Neural Networks*, 113: 54–71.
- Parsons, L.; Haque, E.; and Liu, H. 2004. Subspace clustering for high dimensional data: a review. *Acm sigkdd explorations newsletter*, 6(1): 90–105.
- Qin, Y.; Menara, T.; Oymak, S.; Ching, S.; and Pasqualetti, F. 2022. Non-Stationary Representation Learning in Sequential Linear Bandits. *IEEE Open Journal of Control Systems*.
- Qin, Z.; Cheng, Y.; Zhao, Z.; Chen, Z.; Metzler, D.; and Qin, J. 2020. Multitask mixture of sequential experts for user activity streams. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, 3083–3091.
- Ramesh, R.; and Chaudhari, P. 2021a. Boosting a model zoo for multi-task and continual learning. *arXiv preprint arXiv:2106.03027*.
- Ramesh, R.; and Chaudhari, P. 2021b. Model Zoo: A Growing Brain That Learns Continually. In *International Conference on Learning Representations*.
- Rosenbaum, C.; Klinger, T.; and Riemer, M. 2017. Routing networks: Adaptive selection of non-linear functions for multi-task learning. *arXiv preprint arXiv:1711.01239*.
- Shu, Y.; Kou, Z.; Cao, Z.; Wang, J.; and Long, M. 2021. Zoo-tuning: Adaptive transfer from a zoo of models. In *International Conference on Machine Learning*, 9626–9637. PMLR.

- Shui, C.; Abbasi, M.; Robitaille, L.-É.; Wang, B.; and Gagné, C. 2019. A principled approach for learning task similarity in multitask learning. *arXiv preprint arXiv:1903.09109*.
- Sodhani, S.; Zhang, A.; and Pineau, J. 2021. Multi-task reinforcement learning with context-based representations. In *International Conference on Machine Learning*, 9767–9779. PMLR.
- Strezoski, G.; Noord, N. v.; and Worring, M. 2019. Many task learning with task routing. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 1375–1384.
- Sun, Y.; Narang, A.; Gulluk, I.; Oymak, S.; and Fazel, M. 2021. Towards sample-efficient overparameterized meta-learning. *Advances in Neural Information Processing Systems*, 34: 28156–28168.
- Talagrand, M. 2006. *The generic chaining: upper and lower bounds of stochastic processes*. Springer Science & Business Media.
- Tan, A. Z.; Yu, H.; Cui, L.; and Yang, Q. 2022. Towards personalized federated learning. *IEEE Transactions on Neural Networks and Learning Systems*.
- Thrun, S.; and Pratt, L. 2012. *Learning to learn*. Springer Science & Business Media.
- Tripuraneni, N.; Jin, C.; and Jordan, M. 2021. Provable meta-learning of linear representations. In *International Conference on Machine Learning*, 10434–10443. PMLR.
- Tripuraneni, N.; Jordan, M.; and Jin, C. 2020. On the theory of transfer learning: The importance of task diversity. *Advances in Neural Information Processing Systems*, 33: 7852–7862.
- Vershynin, R. 2010. Introduction to the non-asymptotic analysis of random matrices. *arXiv preprint arXiv:1011.3027*.
- Vershynin, R. 2018. *High-dimensional probability: An introduction with applications in data science*, volume 47. Cambridge university press.
- Vidal, R. 2011. Subspace clustering. *IEEE Signal Processing Magazine*, 28(2): 52–68.
- Vuorio, R.; Sun, S.-H.; Hu, H.; and Lim, J. J. 2019. Multi-modal model-agnostic meta-learning via task-aware modulation. *Advances in Neural Information Processing Systems*, 32.
- Wainwright, M. J. 2019. *High-dimensional statistics: A non-asymptotic viewpoint*, volume 48. Cambridge University Press.
- Wang, J.; Kolar, M.; and Srebro, N. 2016. Distributed multi-task learning with shared representation. *arXiv preprint arXiv:1603.02185*.
- Wortsman, M.; Ramanujan, V.; Liu, R.; Kembhavi, A.; Rastegari, M.; Yosinski, J.; and Farhadi, A. 2020. Supermasks in superposition. *Advances in Neural Information Processing Systems*, 33: 15173–15184.
- Wu, S.; Zhang, H. R.; and Ré, C. 2020. Understanding and improving information transfer in multi-task learning. *arXiv preprint arXiv:2005.00944*.
- Xu, Z.; and Tewari, A. 2021. Representation learning beyond linear prediction functions. *Advances in Neural Information Processing Systems*, 34: 4792–4804.
- Xu, Z.; and Tewari, A. 2022. On the statistical benefits of curriculum learning. In *International Conference on Machine Learning*, 24663–24682. PMLR.
- Yang, J.; Hu, W.; Lee, J. D.; and Du, S. S. 2020. Impact of representation learning in linear bandits. *arXiv preprint arXiv:2010.06531*.
- Yao, H.; Wei, Y.; Huang, J.; and Li, Z. 2019. Hierarchically structured meta-learning. In *International Conference on Machine Learning*, 7045–7054. PMLR.
- Ye, Q.; Zha, J.; and Ren, X. 2022. Eliciting Transferability in Multi-task Learning with Task-level Mixture-of-Experts. *arXiv preprint arXiv:2205.12701*.
- Yu, Y.; Wang, T.; and Samworth, R. J. 2015. A useful variant of the Davis–Kahan theorem for statisticians. *Biometrika*, 102(2): 315–323.
- Zamir, A. R.; Sax, A.; Shen, W.; Guibas, L. J.; Malik, J.; and Savarese, S. 2018. Taskonomy: Disentangling task transfer learning. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 3712–3722.
- Zhang, Y.; and Yang, Q. 2021. A survey on multi-task learning. *IEEE Transactions on Knowledge and Data Engineering*.
- Zhou, F.; Shui, C.; Abbasi, M.; Robitaille, L.-É.; Wang, B.; and Gagné, C. 2020. Task similarity estimation through adversarial multitask neural network. *IEEE Transactions on Neural Networks and Learning Systems*, 32(2): 466–480.
- Zhuang, F.; Qi, Z.; Duan, K.; Xi, D.; Zhu, Y.; Zhu, H.; Xiong, H.; and He, Q. 2020. A comprehensive survey on transfer learning. *Proceedings of the IEEE*, 109(1): 43–76.

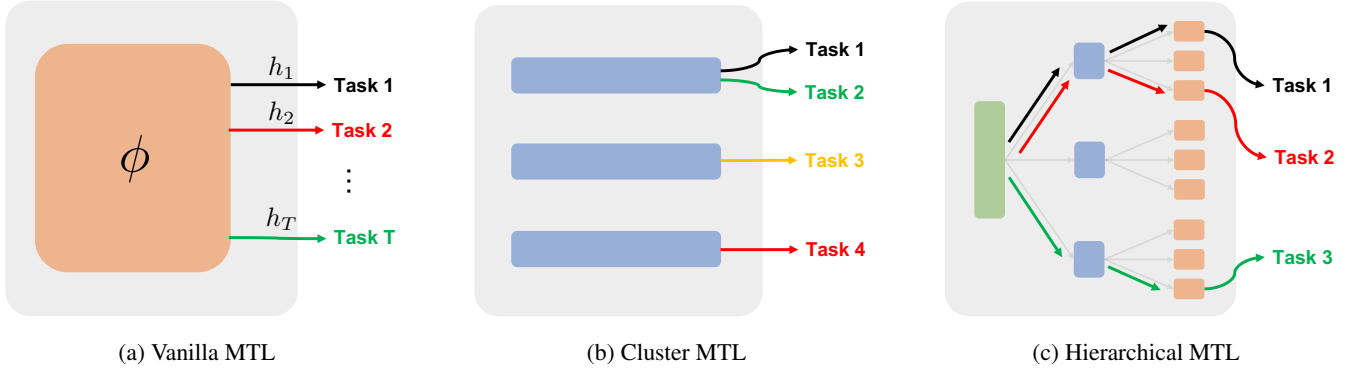


Figure 4: Three specific MTL settings: Vanilla MTL, Cluster MTL and Hierarchical MTL. In vanilla MTL, all the tasks share the same representation  $\phi \in \Phi$ , and each task learns its specific head  $h_t \in \mathcal{H}$ . It corresponds to the setting that  $|\mathcal{A}| = 1$ ,  $L = 1$  and  $K_1 = 1$ . In Cluster MTL, tasks are clustered into groups and different groups are assigned with different and uncorrelated representations. If we assume there are  $K$  clusters, then  $|\mathcal{A}| = K$ ,  $L = 1$  and  $K_1 = K$ . While, Fig. 1b shows the Hierarchical MTL with only two layers, here we present the more general Hierarchical MTL setting. Assume the degree of a hierarchical supernet is  $K$  (In Fig. 4c,  $K = 3$ ), then  $|\mathcal{A}| = K^{L-1}$  where  $L$  is the number of layers in supernet, and  $K_\ell = K^{\ell-1}$ .

## Organization of the Supplementary Material

The supplementary material (SM) is organized as follows.

1. In Appendix A we introduce additional notions used throughout the supplementary material.
2. Appendix B provides our main proofs in Section 3 and introduces two direct corollaries of Theorem 1. We also provide a data-dependent bound in terms of empirical Gaussian complexity (rather than worst-case). In Appendix B.5 we also provide end-to-end transfer learning bound by introducing a proper notion of task diversity.
3. Appendix C provides additional guarantees (Thm 8) for parametric classes via non-data-dependent covering argument. The advantages of Theorem 8 are: (1) Sample complexity has linear dependence on supernet depth  $L$  (rather than exponential), (2) It applies to unbounded loss functions, (3) It is also a supporting result for the proof of Theorem 3&4.
4. Appendix D provides our proofs in Section 4. We also introduce Corollary 4, which is a direct application of Theorem 1. Lemma 7 proves the necessity of our Assumption 5.
5. In Appendix E, we include a short discussion on the challenges of transfer learning: Specifically, we provide a lemma/example that shows that, under the assumptions of Theorem 4, if ground-truth MTL pathways are not known, there are MTL settings for which transfer learning can provably fail. This construction highlights the (combinatorial) challenge of finding the right task clusterings during MTL phase that are actually useful for transfer phase.
6. Appendix F provides further details, algorithms, and results on numerical experiments in Section 5.

## A Useful Definitions

We will start with some useful notions. Let  $\|\cdot\|$  denote the  $\ell_2$ -norm of a vector, and  $[L]$  denote the set  $\{1, 2, \dots, L\}$ . We denote the  $K$  times Cartesian product of a hypothesis set  $\mathcal{Q}$  with itself by  $\mathcal{Q}^K$ . Now assume we have a hypothesis set  $\mathcal{Q} : \mathcal{X} \rightarrow \mathbb{R}^r$  and an input dataset of size  $n$ , defined by  $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ , where  $\mathbf{x}_i \in \mathcal{X}$ . Let  $\{\sigma_{ij}\}_{i \in [n], j \in [r]}$  denote Rademacher variables uniformly and independently taking values in  $\{-1, 1\}$  and  $\{g_{ij}\}_{i \in [n], j \in [r]}$  denote i.i.d. standard random Gaussian variables. Then we can define the empirical and population Rademacher/Gaussian complexities of a hypothesis set  $\mathcal{Q}$  over inputs  $\mathbf{X}$  and data space  $\mathcal{X}$  with sample size  $n$  as

$$\text{Empirical/Population Rademacher complexities: } \hat{\mathcal{R}}_{\mathbf{X}}(\mathcal{Q}) = \mathbb{E}_{\sigma_{ij}} \left[ \sup_{q \in \mathcal{Q}} \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^r \sigma_{ij} q_j(\mathbf{x}_i) \right], \quad \mathcal{R}_n(\mathcal{Q}) = \mathbb{E}_{\mathbf{X}} \left[ \hat{\mathcal{R}}_{\mathbf{X}}(\mathcal{Q}) \right],$$

$$\text{Empirical/Population Gaussian complexities: } \hat{\mathcal{G}}_{\mathbf{X}}(\mathcal{Q}) = \mathbb{E}_{g_{ij}} \left[ \sup_{q \in \mathcal{Q}} \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^r g_{ij} q_j(\mathbf{x}_i) \right], \quad \mathcal{G}_n(\mathcal{Q}) = \mathbb{E}_{\mathbf{X}} \left[ \hat{\mathcal{G}}_{\mathbf{X}}(\mathcal{Q}) \right],$$

where we have  $q \in \mathcal{Q}$  and  $q(\mathbf{x}) = [q_1(\mathbf{x}), \dots, q_r(\mathbf{x})]^\top$ . Note that in vector notation one can also write  $\hat{\mathcal{R}}_{\mathbf{X}}(\mathcal{Q}) = \mathbb{E}_{\sigma_i} \left[ \sup_{q \in \mathcal{Q}} \frac{1}{n} \sum_{i=1}^n \sigma_i^\top q(\mathbf{x}_i) \right]$  and  $\hat{\mathcal{G}}_{\mathbf{X}}(\mathcal{Q}) = \mathbb{E}_{g_i} \left[ \sup_{q \in \mathcal{Q}} \frac{1}{n} \sum_{i=1}^n g_i^\top q(\mathbf{x}_i) \right]$ , where  $\sigma_i$  and  $g_i$  are  $r$ -dimensional with independent Rademacher/Gaussian variables in each entry. Also recall that worst-case versions  $\tilde{\mathcal{R}}_n, \tilde{\mathcal{G}}_n$  are obtained by taking supremum over the input space.

## B Proofs in Section 3

We first introduce some lemmas used throughout this section, then provide the proofs of our mean results.

### B.1 Supporting Lemmas

The following is a seminal contraction lemma due to Talagrand (Talagrand 2006).

**Lemma 2 (Talagrand's Contraction inequality)** *Let  $\varepsilon = (\varepsilon_i)_{i=1}^n$  be i.i.d. random variables with symmetric sign (e.g. Rademacher, standard normal). Let  $(\phi_i)_{i=1}^n$  be  $L$ -Lipschitz functions and  $\mathcal{F}$  be a hypothesis set. We have that*

$$\mathbb{E}_\varepsilon \left[ \sup_{f \in \mathcal{F}} \sum_{i=1}^n \varepsilon_i \phi_i(f(\mathbf{x}_i)) \right] \leq L \mathbb{E}_\varepsilon \left[ \sup_{f \in \mathcal{F}} \sum_{i=1}^n \varepsilon_i f(\mathbf{x}_i) \right].$$

As a corollary of this, we can deduce that adjusted empirical Gaussian complexity  $n\widehat{\mathcal{G}}_{\mathbf{X}}(\mathcal{F})$  is non-decreasing in sample size  $n$ .

**Corollary 1** *Let  $\mathcal{X}$  be a bounded input space and  $\mathcal{F} : \mathcal{X} \rightarrow \mathbb{R}$  be a hypothesis set. Let  $\mathbf{X}_m$  be a dataset of size  $m$  and  $\mathbf{X}_n = (\mathbf{x}_i)_{i=1}^n$  be a dataset of size  $n$  that contains  $\mathbf{X}_m$ . We have that*

$$m\widehat{\mathcal{G}}_{\mathbf{X}_m}(\mathcal{F}) \leq n\widehat{\mathcal{G}}_{\mathbf{X}_n}(\mathcal{F}).$$

We note that, when  $\mathcal{F} : \mathcal{X} \rightarrow \mathbb{R}^p$  is vector valued and we apply  $p \times n$   $L$ -Lipschitz functions  $\phi_{ij}$ , the identical results (Lemmas 2 and Corollary 1) follow from Sudakov-Fernique inequality under Gaussian  $\varepsilon \in \mathbb{R}^{n \times p}$  (e.g. Exercise 7.2.13 of (Vershynin 2018)).

This also implies usual (distributional) and worst-case Gaussian complexities are also non-decreasing.

**Proof** Let  $(\phi_i)_{i=1}^n$  be functions that are identity for  $i \leq m$  and zero for  $i > m$ . Observe that

$$m\widehat{\mathcal{G}}_{\mathbf{X}_m} = \mathbb{E}_\varepsilon \left[ \sup_{f \in \mathcal{F}} \sum_{i=1}^m \varepsilon_i f(\mathbf{x}_i) \right] = \mathbb{E}_\varepsilon \left[ \sup_{f \in \mathcal{F}} \sum_{i=1}^n \varepsilon_i \phi_i(f(\mathbf{x}_i)) \right] \leq n\widehat{\mathcal{G}}_{\mathbf{X}_n}. \quad \blacksquare$$

The following lemma shows that adjusted worst-case Gaussian complexity  $\sqrt{n}\widetilde{\mathcal{G}}_{\mathbf{X}}(\mathcal{F})$  is essentially non-decreasing in sample size  $n$ .

**Lemma 3 (Worst-case Gaussian Complexity over Input Space and Sample Size)** *For any bounded input space  $\mathcal{X}$  and hypothesis set  $\mathcal{F}$ , we have that*

$$\sup_{1 \leq m \leq n} \sqrt{m}\widetilde{\mathcal{G}}_m(\mathcal{F}) \leq \sqrt{2n}\widetilde{\mathcal{G}}_n(\mathcal{F}).$$

**Proof** First suppose  $n/2 \leq m \leq n$ . In this case, from Corollary 1, we know that  $m\widetilde{\mathcal{G}}_m(\mathcal{F}) \leq n\widetilde{\mathcal{G}}_n(\mathcal{F}) \implies \sqrt{m}\widetilde{\mathcal{G}}_m(\mathcal{F}) \leq \frac{\sqrt{2m}}{\sqrt{n}}\widetilde{\mathcal{G}}_m(\mathcal{F}) \leq \sqrt{2n}\widetilde{\mathcal{G}}_n(\mathcal{F})$ . What remains is the scenario  $m < n/2$ . To do this, we will show monotonicity under doubling  $\sqrt{m}\widetilde{\mathcal{G}}_m(\mathcal{F}) \leq \sqrt{2m}\widetilde{\mathcal{G}}_{2m}(\mathcal{F})$ . If this holds, then you can double  $m$  until a point  $n/2 \leq m \leq n$  and apply the first bound.

Consider worst-case dataset for  $\widetilde{\mathcal{G}}_m$  defined as

$$\mathbf{Y} = \arg \max_{\mathbf{X} \in \mathcal{X}^m} \widehat{\mathcal{G}}_{\mathbf{X}}(\mathcal{F}).$$

Let  $\mathbf{Y}'$  be a dataset of size  $2m$  that repeats the elements of  $\mathbf{Y}$  twice so that  $\mathbf{y}'_{m+i} = \mathbf{y}'_i = \mathbf{y}_i$ . Here, we consider hypothesis set  $\mathcal{F} : \mathcal{X} \rightarrow \mathbb{R}^p$ , and then  $f(\mathbf{y}_i) = [f_1(\mathbf{y}_i), \dots, f_p(\mathbf{y}_i)]^\top$ . Also let  $\varepsilon \in \mathbb{R}^{m \times p}$ ,  $\varepsilon' \in \mathbb{R}^{2m \times p}$  where  $\varepsilon'_i \sim \mathcal{N}(0, I_p)$ ,  $i \in [2m]$  and  $\varepsilon_i = \frac{\varepsilon'_i + \varepsilon'_{m+i}}{\sqrt{2}} \sim \mathcal{N}(0, I_p)$ . We have that

$$\begin{aligned} 2m\widetilde{\mathcal{G}}_{2m}(\mathcal{F}) &\geq \mathbb{E}_\varepsilon \left[ \sup_{f \in \mathcal{F}} \sum_{i=1}^{2m} \sum_{j=1}^p \varepsilon'_{ij} f_j(\mathbf{y}'_i) \right] \\ &\geq \mathbb{E}_\varepsilon \left[ \sup_{f \in \mathcal{F}} \sum_{i=1}^m \sum_{j=1}^p \varepsilon'_{i,j} f_j(\mathbf{y}'_i) + \varepsilon'_{(m+i),j} f_j(\mathbf{y}'_{m+i}) \right] \\ &\geq \sqrt{2} \mathbb{E}_\varepsilon \left[ \sup_{f \in \mathcal{F}} \sum_{i=1}^m \sum_{j=1}^p \varepsilon_{ij} f_j(\mathbf{y}_i) \right] \\ &= \sqrt{2}m\widetilde{\mathcal{G}}_m. \end{aligned}$$

Dividing both sides by  $\sqrt{2m}$ , we conclude with the claim  $\sqrt{m}\widetilde{\mathcal{G}}_m(\mathcal{F}) \leq \sqrt{2m}\widetilde{\mathcal{G}}_{2m}(\mathcal{F})$ . \(\blacksquare\)

The following is a model selection argument shows that  $\widetilde{\mathcal{G}}(\Phi)$  can be replaced with  $\widetilde{\mathcal{G}}(\Phi_{\text{used}})$ .

**Lemma 4 (Only utilized supernet matters)** *Observe that  $T$  tasks can choose from up to  $|\mathcal{A}|^T$  supernets in total. Let  $\Phi_{all} = (\Phi_i)_{i=1}^H$  with  $H \leq |\mathcal{A}|^T$  be the set of unique supernets (since two supernets that choose same number of modules per layer are identical architectures). Suppose the outcome of empirical risk minimization ( $M^2TL$ ) obeys  $\hat{\phi} \in \Phi_{used} \in \Phi_{all}$ . Let  $\hat{K}_\ell$  be the number of (used) modules in  $\Phi_{used}$ . With probability  $1 - \delta$ , we have that*

$$\mathcal{L}_{\mathcal{D}}(\hat{\mathbf{f}}) - \hat{\mathcal{L}}_{\mathcal{S}_{all}}(\hat{\mathbf{f}}) \lesssim \tilde{\mathcal{G}}_N(\mathcal{H}) + \sum_{\ell=1}^L \sqrt{\hat{K}_\ell} \tilde{\mathcal{G}}_{NT}(\Psi_\ell) + \sqrt{\frac{\log |\mathcal{A}|}{N} + \frac{\log(2/\delta)}{NT}}, \quad (5)$$

$$\mathcal{R}_{M^2TL}(\hat{\mathbf{f}}) := \mathcal{L}_{\mathcal{D}}(\hat{\mathbf{f}}) - \mathcal{L}_{\mathcal{D}}^*(\hat{\mathbf{f}}) \lesssim \tilde{\mathcal{G}}_N(\mathcal{H}) + \sum_{\ell=1}^L \sqrt{\hat{K}_\ell} \tilde{\mathcal{G}}_{NT}(\Psi_\ell) + \sqrt{\frac{\log |\mathcal{A}|}{N} + \frac{\log(2/\delta)}{NT}}. \quad (6)$$

**Proof** Let  $\mathcal{L}_{\Phi'}$ ,  $\hat{\mathcal{L}}_{\Phi'}$  be the population and empirical risks we achieve when we run the ( $M^2TL$ ) problem over  $\Phi' \in \Phi_{all}$  rather than  $\Phi$ . Additionally, let  $K_\ell(\Phi')$  denote the number of modules in the  $\ell$ th layer of the architecture  $\Phi'$ . Given  $\Phi'$ , also define  $\mathcal{C}_N(\Phi', \delta)$  to be the excess risk bound one obtains via (9) ((9) in Theorem 5 is obtained without using Lemma 4), that is,

$$\mathcal{C}_N(\Phi', \delta) = \tilde{\mathcal{G}}_N(\mathcal{H}) + \sum_{\ell=1}^L \sqrt{K_\ell(\Phi')} \tilde{\mathcal{G}}_{NT}(\Psi_\ell) + \sqrt{\frac{\log |\mathcal{A}|}{N} + \frac{\log(2/\delta)}{NT}}.$$

To proceed, applying (9) over  $\Phi' \in \Phi_{all}$  and union bounding over all  $H \leq |\mathcal{A}|^T$ , with probability at least  $1 - \delta$ , we find that, all  $\Phi' \in \Phi_{all}$  obeys

$$|\hat{\mathcal{L}}_{\Phi'}(\hat{\mathbf{f}}) - \mathcal{L}_{\Phi'}(\hat{\mathbf{f}})| \lesssim \mathcal{C}_N(\Phi', \delta/H).$$

Fortunately,  $\mathcal{C}_N(\Phi', \delta/H) \lesssim \mathcal{C}_N(\Phi', \delta)$  since the latter already includes a  $\sqrt{\frac{\log |\mathcal{A}|}{N}}$  term. Using this union bound, optimality of  $\hat{\phi} \in \Phi_{used}$  (and that of the associated  $\hat{\mathbf{f}} \in \mathcal{F}_{used}$ ), and using  $\hat{\mathcal{L}}_{\Phi_{used}}(\hat{\mathbf{f}}) = \hat{\mathcal{L}}_{\Phi}(\hat{\mathbf{f}}) = \hat{\mathcal{L}}_{\mathcal{S}_{all}}(\hat{\mathbf{f}})$ , we find that

$$\mathcal{L}_{\Phi_{used}}(\hat{\mathbf{f}}) \leq \hat{\mathcal{L}}_{\Phi_{used}}(\hat{\mathbf{f}}) + \mathcal{O}(\mathcal{C}_N(\Phi_{used}, \delta)) \quad (7)$$

$$\leq \hat{\mathcal{L}}_{\mathcal{S}_{all}}(\hat{\mathbf{f}}) + \mathcal{O}(\mathcal{C}_N(\Phi_{used}, \delta)). \quad (8)$$

The last line establishes Inequality (5). To conclude with the second inequality, we control the excess risk error by observing test risk upper bounds the training risk. Namely, let  $\mathbf{f}_\star \in \mathcal{F}$  be the population minima. First, with  $1 - \delta$  probability, for this singleton hypothesis, we have that

$$|\mathcal{L}_{\mathcal{D}}(\mathbf{f}_\star) - \hat{\mathcal{L}}_{\mathcal{S}_{all}}(\mathbf{f}_\star)| \leq \sqrt{\frac{\log(2/\delta)}{NT}}.$$

Second, we can write

$$\hat{\mathcal{L}}_{\mathcal{S}_{all}}(\hat{\mathbf{f}}) \leq \hat{\mathcal{L}}_{\mathcal{S}_{all}}(\mathbf{f}_\star) \leq \mathcal{L}_{\mathcal{D}}(\mathbf{f}_\star) + \sqrt{\frac{\log(2/\delta)}{NT}}.$$

Combining this with (8), we establish the guarantee against the ground-truth optima  $\mathbf{f}_\star$

$$\mathcal{L}_{\Phi_{used}}(\hat{\mathbf{f}}) - \left[ \mathcal{L}_{\mathcal{D}}(\mathbf{f}_\star) + \sqrt{\frac{\log(2/\delta)}{NT}} \right] \leq \mathcal{L}_{\Phi_{used}}(\hat{\mathbf{f}}) - \hat{\mathcal{L}}_{\mathcal{S}_{all}}(\hat{\mathbf{f}}) \leq \mathcal{O}(\mathcal{C}_N(\Phi_{used}, \delta)),$$

which establishes the claim (6) after subsuming  $\sqrt{\frac{\log(2/\delta)}{NT}}$  within  $\mathcal{C}_N(\Phi_{used}, \delta)$ . ■

## B.2 Proof of Theorem 1

Let us define the covering number of a hypothesis as well as natural data-dependent Euclidean distance for ease of reference in the subsequent discussion (see (Wainwright 2019)).

**Definition 3 (Covering number)** *Let  $\mathcal{Q} : \mathcal{X} \rightarrow \mathbb{R}^r$  be a family of functions. Given  $q, q' \in \mathcal{Q}$ , and a distance metric  $d(q, q') \geq 0$ , an  $\varepsilon$ -cover of set  $\mathcal{Q}$  with respect to  $d(\cdot, \cdot)$  is a set  $\{q^1, q^2, \dots, q^N\} \subset \mathcal{Q}$  such that for any  $q \in \mathcal{Q}$ , there exists some  $i \in [N]$  such that  $d(q, q^i) \leq \varepsilon$ . The  $\varepsilon$ -covering number  $\mathcal{N}(\varepsilon; \mathcal{Q}, d)$  is defined to be the cardinality of the smallest  $\varepsilon$ -cover.*

**Definition 4 (Data-dependent distance metric  $\rho$ )** *Let  $\mathcal{Q} : \mathcal{X} \rightarrow \mathbb{R}^r$  be a family of functions. Given  $q, q' \in \mathcal{Q}$  and an input dataset  $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  with  $\mathbf{x}_i \in \mathcal{X}$ , we define the dataset-dependent Euclidean distance by  $\rho_{\mathbf{X}}(q, q') := \sqrt{\frac{1}{n} \sum_{i \in [n], j \in [r]} (q_j(\mathbf{x}_i) - q'_j(\mathbf{x}_i))^2} = \sqrt{\frac{1}{n} \sum_{i \in [n]} \|q(\mathbf{x}_i) - q'(\mathbf{x}_i)\|^2}$ , where  $q(\mathbf{x}) = [q_1(\mathbf{x}), \dots, q_r(\mathbf{x})]^\top$ .*

Now we are ready to prove our main theorem which incorporates additional dependencies that were omitted from the original statement.

**Theorem 5 (Theorem 1 restated)** Suppose Assumptions 1&2 hold. Let  $\hat{\mathbf{f}}$  be the empirical solution of (M<sup>2</sup>TL). Let  $D_{\mathcal{X}} = \sup_{\mathbf{x} \in \mathcal{X}, h \in \mathcal{H}, \phi \in \Phi, \alpha \in \mathcal{A}} |h \circ \phi_{\alpha}(\mathbf{x})| < \infty$ , and set  $\Gamma^{\dagger} = \sum_{\ell=0}^L \Gamma^{\ell}$ . Then, with probability at least  $1 - \delta$ , the excess test risk in (2) obeys

$$\mathcal{R}_{M^2TL}(\hat{\mathbf{f}}) \leq 768\Gamma \left( \frac{D_{\mathcal{X}}}{NT} + D_{\mathcal{X}} \sqrt{\frac{\log |\mathcal{A}|}{N}} + \Gamma^{\dagger} \log NT \left( \tilde{\mathcal{G}}_N(\mathcal{H}) + \sum_{\ell=1}^L \sqrt{K_{\ell}} \tilde{\mathcal{G}}_{NT}(\Psi_{\ell}) \right) \right) + 2\sqrt{\frac{\log \frac{2}{\delta}}{NT}}. \quad (9)$$

Here, the input spaces for  $\mathcal{H}$  and  $\Psi_{\ell}$  are  $\mathcal{X}_{\mathcal{H}} = \Psi_L \circ \dots \circ \Psi_1 \circ \mathcal{X}$ ,  $\mathcal{X}_{\Psi_{\ell}} = \Psi_{\ell-1} \circ \dots \circ \Psi_1 \circ \mathcal{X}$  for  $\ell > 1$ , and  $\mathcal{X}_{\Psi_1} = \mathcal{X}$ . The above is our general results, which we do not focus on the actual modules used in  $\hat{\mathbf{f}}$ . Now let  $\hat{K}$  be the number of modules utilized by  $\hat{\mathbf{f}}$ , then with probability at least  $1 - \delta$ , we can obtain

$$\mathcal{R}_{M^2TL}(\hat{\mathbf{f}}) \lesssim \tilde{\mathcal{G}}_N(\mathcal{H}) + \sum_{\ell=1}^L \sqrt{\hat{K}_{\ell}} \tilde{\mathcal{G}}_{NT}(\Psi_{\ell}) + \sqrt{\frac{\log |\mathcal{A}|}{N} + \frac{\log(2/\delta)}{NT}}. \quad (10)$$

Here,  $\lesssim$  suppresses dependencies on  $\log NT$ ,  $\Gamma^{\dagger}$  and  $D_{\mathcal{X}}$ .

**Remark.** While this result is stated with worst-case Gaussian complexity, the line (20) states our result in terms of empirical Gaussian complexity which is always a lower bound and is in terms of the training dataset. However, (20) is more convoluted and involves worst-case hypothesis being applied to the training data. The latter arises from the fact that it is difficult to track the evolution of features across arbitrary pathways and hierarchical layers.

**Proof** To start with, let us recap some notations. Assume we have  $T$  tasks each with  $N$  training samples i.i.d. drawn from  $(\mathcal{D}_t)_{t=1}^T$  respectively, and let  $\bar{\mathcal{D}} = \{\mathcal{D}_t\}_{t=1}^T$ . Denote the training dataset and inputs of  $t_{\text{th}}$  task by  $\mathcal{S}_t = \{(\mathbf{x}_{ti}, y_{ti})\}_{i=1}^N$  and  $\mathbf{X}_t = \{\mathbf{x}_{ti}\}_{i=1}^N$ , and define the union by  $\mathcal{S}_{\text{all}} = \bigcup_{t=1}^T \mathcal{S}_t$  and  $\mathbf{X} = \bigcup_{t=1}^T \mathbf{X}_t$ . Let  $\mathbf{h} = [h_1, \dots, h_T] \in \mathcal{H}^T$ ,  $\boldsymbol{\alpha} = [\alpha_1, \dots, \alpha_T] \in \mathcal{A}^T$ ,  $\boldsymbol{\psi}_{\ell} = [\psi_{\ell}^1, \dots, \psi_{\ell}^{K_{\ell}}] \in \Psi_{\ell}^{K_{\ell}}$ ,  $\ell \in [L]$ , and  $\boldsymbol{\phi} = [\psi_1, \dots, \psi_L] \in \Phi = \Psi_1^{K_1} \times \dots \times \Psi_L^{K_L}$ .  $\hat{\mathbf{f}} := (\hat{\mathbf{h}}, \hat{\boldsymbol{\alpha}}, \hat{\boldsymbol{\phi}})$  is the empirical solution of (M<sup>2</sup>TL) and  $\mathbf{f}^* := (\mathbf{h}^*, \boldsymbol{\alpha}^*, \boldsymbol{\phi}^*)$  is the population solution of (M<sup>2</sup>TL) when each task has infinite i.i.d training samples ( $N = \infty$ ). Let  $\mathcal{F}$  denote the hypothesis set of functions  $\mathbf{f}$ . Since multitask problem is task-aware, that is, the task identification of each data is given during training and test, we can rewrite samples in  $\mathcal{S}_t$  as  $\{(\mathbf{x}_i, y_i, t_i \equiv t)\}_{i=1+(t-1)N}^{tN}$  and the overall multitask training dataset can be seen as  $\mathcal{S}_{\text{all}} = \{(\mathbf{x}_i, y_i, t_i)\}_{i=1}^{NT}$ . Letting  $\mathbf{f}(\mathbf{x}, t) = f_t(\mathbf{x}) = h_t \circ \phi_{\alpha_t}(\mathbf{x})$ , the loss functions can be rewritten by  $\hat{\mathcal{L}}_{\mathcal{S}_{\text{all}}}(\mathbf{f}) = \frac{1}{NT} \sum_{i=1}^{NT} \ell(\mathbf{f}(\mathbf{x}_i, t_i), y_i)$  and  $\mathcal{L}_{\bar{\mathcal{D}}}(\mathbf{f}) = \mathbb{E}[\hat{\mathcal{L}}_{\mathcal{S}_{\text{all}}}(\mathbf{f})]$ . In the following, we drop the subscript  $\bar{\mathcal{D}}$  and  $\mathcal{S}_{\text{all}}$  for cleaner notations. Then we have

$$\underbrace{\mathcal{L}(\hat{\mathbf{f}}) - \mathcal{L}(\mathbf{f}^*)}_{\mathcal{R}_{M^2TL}(\hat{\mathbf{f}})} = \underbrace{\mathcal{L}(\hat{\mathbf{f}}) - \hat{\mathcal{L}}(\hat{\mathbf{f}})}_a + \underbrace{\hat{\mathcal{L}}(\hat{\mathbf{f}}) - \hat{\mathcal{L}}(\mathbf{f}^*)}_b + \underbrace{\hat{\mathcal{L}}(\mathbf{f}^*) - \mathcal{L}(\mathbf{f}^*)}_c, \quad (11)$$

where  $b \leq 0$  because of the fact that  $\hat{\mathbf{f}}$  is the empirical risk minimizer of  $\hat{\mathcal{L}}(\mathbf{f})$ . Then, following the proof of Theorem 3.3 of (Mohri, Rostamizadeh, and Talwalkar 2018), we make two observations: 1) Their Equation (3.8) in the proof still holds when we restrict  $N$  i.i.d samples in each task instead of  $NT$  i.i.d. samples over distribution  $\bar{\mathcal{D}}$ . Therefore, the symmetrization augment does not change, and this theorem holds under our setting. 2) The identical results hold for any function set mapping to  $[-1, 1]$ . In this work, based on these two observations, following Assumption 2 and Theorem 11.3 in (Mohri, Rostamizadeh, and Talwalkar 2018), we have that with probability at least  $1 - \delta/2$ ,  $a, c \leq 2\Gamma \mathcal{R}_{NT}(\mathcal{F}) + \sqrt{\frac{\log(2/\delta)}{2NT}}$ . Therefore, we can conclude that with probability at least  $1 - \delta$ ,

$$\mathcal{R}_{M^2TL}(\hat{\mathbf{f}}) \leq 4\Gamma \mathcal{R}_{NT}(\mathcal{F}) + \sqrt{\frac{2 \log \frac{2}{\delta}}{NT}}, \quad (12)$$

$$\text{and similarly, } \mathcal{R}_{M^2TL}(\hat{\mathbf{f}}) \leq 4\Gamma \hat{\mathcal{R}}_{\mathbf{X}}(\mathcal{F}) + 3\sqrt{\frac{2 \log \frac{4}{\delta}}{NT}}, \quad (13)$$

where  $\hat{\mathcal{R}}_{\mathbf{X}}(\mathcal{F})$  is the empirical complexity with respect to the inputs  $\mathbf{X}$  and  $\mathcal{R}_{NT}(\mathcal{F})$  is the Rademacher complexity with respect to the sample size  $NT$ . Exercise 5.5 in (Wainwright 2019) shows that Rademacher complexity can be bounded in terms of Gaussian complexity, that is  $\hat{\mathcal{R}}_{\mathbf{X}}(\mathcal{F}) \leq \sqrt{\frac{\pi}{2}} \hat{\mathcal{G}}_{\mathbf{X}}(\mathcal{F})$  and  $\mathcal{R}_{NT}(\mathcal{F}) \leq \sqrt{\frac{\pi}{2}} \mathcal{G}_{NT}(\mathcal{F})$ . Combining them together, we have that with probability at least  $1 - \delta$ ,

$$\mathcal{R}_{M^2TL}(\hat{\mathbf{f}}) \leq 6\Gamma \mathcal{G}_{NT}(\mathcal{F}) + 2\sqrt{\frac{\log \frac{2}{\delta}}{NT}}, \quad \text{and} \quad \mathcal{R}_{M^2TL}(\hat{\mathbf{f}}) \leq 6\Gamma \hat{\mathcal{G}}_{\mathbf{X}}(\mathcal{F}) + 6\sqrt{\frac{\log \frac{4}{\delta}}{NT}}. \quad (14)$$

In what follows, we will move to Gaussian complexity instead. Now, it remains to decompose the Gaussian complexity of a set of composition functions  $\mathcal{F}$  into basic function sets  $\mathcal{H}$ ,  $\mathcal{A}$  and  $\{\Psi_\ell\}_{\ell=1}^L$ . We will first bound the empirical Gaussian complexity with respect to any training inputs  $\mathbf{X}$ , which turns to be worst-case Gaussian complexity defined in Definition 1. Then, population complexity is simply bounded by the worst-case Gaussian complexity.

Inspired by (Tripuraneni, Jordan, and Jin 2020), we use the Dudley's entropy integral bound showed in (Wainwright 2019) (Theorem 5.22) to derive the upper bound. Define  $Z_{\mathbf{f}} := \frac{1}{\sqrt{NT}} \sum_{i=1}^{NT} g_i \mathbf{f}(\mathbf{x}_i, t_i)$  where  $\mathbf{f} \in \mathcal{F}$  and  $g_i$ s are standard random Gaussian variables. Since  $Z_{\mathbf{f}}$  has zero-mean, we have  $\widehat{\mathcal{G}}_{\mathbf{X}}(\mathcal{F}) = \frac{1}{\sqrt{NT}} \mathbb{E}_{\mathbf{g}}[\sup_{\mathbf{f} \in \mathcal{F}} Z_{\mathbf{f}}] \leq \frac{1}{\sqrt{NT}} \mathbb{E}_{\mathbf{g}}[\sup_{\mathbf{f}, \mathbf{f}' \in \mathcal{F}} (Z_{\mathbf{f}} - Z_{\mathbf{f}'})]$ . Following Definition 4, let  $\rho_{\mathbf{X}}(\mathbf{f}, \mathbf{f}') = \sqrt{\frac{1}{NT} \sum_{i=1}^{NT} (\mathbf{f}(\mathbf{x}_i, t_i) - \mathbf{f}'(\mathbf{x}_i, t_i))^2}$ . Define  $D_{\mathbf{X}} = \sup_{\mathbf{f}, \mathbf{f}' \in \mathcal{F}} \rho_{\mathbf{X}}(\mathbf{f}, \mathbf{f}') \leq 2D_{\mathcal{X}}$ . Following Theorem 5.22 in (Wainwright 2019), we have that for any  $\varepsilon \in [0, D_{\mathbf{X}}]$ ,

$$\mathbb{E}_{\mathbf{g}} \left[ \sup_{\mathbf{f}, \mathbf{f}' \in \mathcal{F}} (Z_{\mathbf{f}} - Z_{\mathbf{f}'}) \right] \leq 2 \mathbb{E}_{\mathbf{g}} \left[ \sup_{\substack{\mathbf{f}, \mathbf{f}' \in \mathcal{F} \\ \rho_{\mathbf{X}}(\mathbf{f}, \mathbf{f}') \leq \varepsilon}} (Z_{\mathbf{f}} - Z_{\mathbf{f}'}) \right] + 32 \int_{\varepsilon/4}^{D_{\mathbf{X}}} \sqrt{\log \mathcal{N}(u; \mathcal{F}, \rho_{\mathbf{X}})} du, \quad (15)$$

where  $\mathcal{N}(u; \mathcal{F}, \rho_{\mathbf{X}})$  is the  $u$ -covering number of function set  $\mathcal{F}$  with respect to metric  $\rho_{\mathbf{X}}(\cdot, \cdot)$  following Definition 3.

The first term in the right hand side above is easy to bound. As shown in proof of Theorem 7 in (Tripuraneni, Jordan, and Jin 2020), we have  $\mathbb{E}_{\mathbf{g}}[\sup_{\rho_{\mathbf{X}}(\mathbf{f}, \mathbf{f}') \leq \varepsilon} (Z_{\mathbf{f}} - Z_{\mathbf{f}'})] \leq \mathbb{E}_{\mathbf{g}}[\sup_{\|\mathbf{v}\|_2 \leq \varepsilon} \mathbf{g}^\top \mathbf{v}] \leq \mathbb{E}_{\mathbf{g}}[\sup_{\|\mathbf{v}\|_2 \leq \varepsilon} \|\mathbf{g}\|_2 \|\mathbf{v}\|_2] = \sqrt{NT} \varepsilon$ . Next, it remains to bound the integral term. Here, since  $\mathbf{f} \in \mathcal{F}$  is a sophisticated function composed with  $\psi_\ell^k \in \Psi_\ell, \alpha_t \in \mathcal{A}$  and  $h_t \in \mathcal{H}$ , its covering number is not well-defined. Hence, instead, we relate the cover of  $\mathcal{F}$  to the covers of basic function sets,  $\Psi_\ell, \mathcal{A}$  and  $\mathcal{H}$ . To this end, we need to decompose the distance metric  $\rho_{\mathbf{X}}$  into distances over basic sets. Since  $\mathcal{A}$  is a discrete set with cardinality  $|\mathcal{A}|$ . Let  $\mathcal{F}^\alpha \subset \mathcal{F}$  be the function set given pathways of all tasks  $\alpha$ . Then we have  $\log \mathcal{N}(u; \mathcal{F}, \rho_{\mathbf{X}}) \leq T \log |\mathcal{A}| + \max_{\alpha \in \mathcal{A}^T} \log \mathcal{N}(u; \mathcal{F}^\alpha, \rho_{\mathbf{X}})$ . For any  $\mathbf{f}, \mathbf{f}' \in \mathcal{F}^\alpha$ , we have

$$\begin{aligned} \rho_{\mathbf{X}}(\mathbf{f}, \mathbf{f}') &= \sqrt{\frac{1}{NT} \sum_{i=1}^{NT} (\mathbf{f}(\mathbf{x}_i, t_i) - \mathbf{f}'(\mathbf{x}_i, t_i))^2} = \sqrt{\frac{1}{NT} \sum_{t=1}^T \sum_{i=1}^N (h_t \circ \phi_{\alpha_t}(\mathbf{x}_{ti}) - h'_t \circ \phi'_{\alpha_t}(\mathbf{x}_{ti}))^2} \\ &\leq \underbrace{\sqrt{\frac{1}{NT} \sum_{t=1}^T \sum_{i=1}^N (h_t \circ \phi_{\alpha_t}(\mathbf{x}_{ti}) - h'_t \circ \phi_{\alpha_t}(\mathbf{x}_{ti}))^2}}_d + \underbrace{\sqrt{\frac{1}{NT} \sum_{t=1}^T \sum_{i=1}^N (h'_t \circ \phi_{\alpha_t}(\mathbf{x}_{ti}) - h'_t \circ \phi'_{\alpha_t}(\mathbf{x}_{ti}))^2}}_e. \end{aligned}$$

To proceed, let us introduce some notations. For any function  $\phi$  with inputs  $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ , define output set w.r.t. the inputs  $\mathbf{X}$  by  $\phi(\mathbf{X}) = \{\phi(\mathbf{x}_1), \dots, \phi(\mathbf{x}_n)\}$ . In the multipath setting, since different tasks have different pathways, different modules are chosen by different set of tasks. Given  $\alpha$ , the task clustering methods in different layers are determined. Let  $\mathcal{I}_\ell^k$  denote the union of task IDs who select  $(\ell, k)$ 'th module, and  $\mathcal{I}_\ell^k, \ell \in [K_\ell]$  are disjoint sets satisfying  $\bigcup_{k=1}^{K_\ell} \mathcal{I}_\ell^k = [T]$ . What's more, let  $\mathbf{Z}_\ell^k$  denote the latent inputs of  $(\ell, k)$ 'th module, where we have

$$\mathbf{Z}_\ell^k = \bigcup_{t \in \mathcal{I}_\ell^k} \psi_{\ell-1}^{\alpha_t} \cdots \circ \psi_1^{\alpha_t}(\mathbf{X}_t), \quad 1 < \ell \leq L, \quad (16)$$

and  $\mathbf{Z}_1^k = \bigcup_{t \in \mathcal{I}_1^k} \mathbf{X}_t$ . In short,  $(\ell, k)$ 'th module (whose function is  $\psi_\ell^k$ ) is utilized by tasks  $\mathcal{I}_\ell^k$  with latent inputs  $\mathbf{Z}_\ell^k$ . The inputs of heads are

$$\mathbf{Z}_{\mathcal{H}}^t = \psi_L^{\alpha_t} \cdots \circ \psi_1^{\alpha_t}(\mathbf{X}_t) = \phi_{\alpha_t}(\mathbf{X}_t), \quad \forall t \in [T].$$



Then we can obtain that

$$\begin{aligned}
(d) &= \sqrt{\frac{1}{T} \sum_{t=1}^T \frac{1}{N} \sum_{i=1}^N (h_t \circ \phi_{\alpha_t}(\mathbf{x}_{ti}) - h'_t \circ \phi_{\alpha_t}(\mathbf{x}_{ti}))^2} \leq \sqrt{\frac{1}{T} \sum_{t=1}^T \rho_{\mathcal{Z}_{\mathcal{H}}^t}^2(h_t, h'_t)}, \\
(e) &\leq \Gamma \sqrt{\frac{1}{NT} \sum_{t=1}^T \sum_{i=1}^N \|\phi_{\alpha_t}(\mathbf{x}_{ti}) - \phi'_{\alpha_t}(\mathbf{x}_{ti})\|^2} \\
&\leq \Gamma \sum_{\ell=1}^L \Gamma^{L-\ell} \sqrt{\frac{1}{K_\ell} \sum_{k=1}^{K_\ell} \frac{1}{|\mathcal{Z}_\ell^k|} \sum_{\mathbf{z}_i \in \mathcal{Z}_\ell^k} \|\psi_\ell^k(\mathbf{z}_i) - \psi'^k_\ell(\mathbf{z}_i)\|^2} \\
&\leq \sum_{\ell=1}^L \Gamma^{L-\ell+1} \sqrt{\frac{1}{K_\ell} \sum_{k=1}^{K_\ell} \rho_{\mathcal{Z}_\ell^k}^2(\psi_\ell^k, \psi'^k_\ell)}.
\end{aligned}$$

Here  $|\mathcal{Z}_\ell^k| = |\mathcal{I}_\ell^k|N$  is the number of samples used in training  $(\ell, k)$ 'th module. The result follows the fact that all functions  $h \in \mathcal{H}, \psi_\ell^k \in \Psi_\ell, \ell \in [L], k \in [K_\ell]$  are  $\Gamma$ -Lipschitz, and it also applies an implicit chain rule for composition Lipschitz functions. Now, we decompose distance  $(d)$  into distances of each head function  $h_t, t \in [T]$ , with inputs  $\mathcal{Z}_{\mathcal{H}}^t$ , and decompose distance  $(e)$ , which captures the distance of composition functions  $\phi$  and  $\phi'$ , into distances of module functions  $\psi_\ell^k, \psi'^k_\ell, \ell \in [L], k \in [K_\ell]$ , w.r.t. inputs of  $\psi_\ell^k, \mathcal{Z}_\ell^k$ . Combining them together and assuming  $\rho_{\mathcal{Z}_{\mathcal{H}}^t}(h_t, h'_t) \leq \varepsilon'$  and  $\rho_{\mathcal{Z}_\ell^k}(\psi_\ell^k, \psi'^k_\ell) \leq \varepsilon'$  for all  $t \in [T], \ell \in [L]$  and  $k \in [K_\ell]$ , we can obtain

$$\rho_{\mathbf{X}}(\mathbf{f}, \mathbf{f}') \leq \sqrt{\frac{1}{T} \sum_{t=1}^T \rho_{\mathcal{Z}_{\mathcal{H}}^t}^2(h_t, h'_t)} + \sum_{\ell=1}^L \Gamma^{L-\ell+1} \sqrt{\frac{1}{K_\ell} \sum_{k=1}^{K_\ell} \rho_{\mathcal{Z}_\ell^k}^2(\psi_\ell^k, \psi'^k_\ell)} \leq \left(1 + \sum_{\ell=1}^L \Gamma^{L-\ell+1}\right) \varepsilon' := \Gamma^\dagger \varepsilon'.$$

It shows that given pathway assignments  $\alpha$  and inputs  $\mathbf{X}$ ,  $\varepsilon'$ -covers of all heads and modules result in  $(\Gamma^\dagger \varepsilon)$ -cover of  $\mathcal{F}^\alpha$ . Recalling that  $\log \mathcal{N}(u; \mathcal{F}, \rho_{\mathbf{X}}) \leq T \log |\mathcal{A}| + \max_{\alpha \in \mathcal{A}^T} \log \mathcal{N}(u; \mathcal{F}^\alpha, \rho_{\mathbf{X}})$ , we have

$$\log \mathcal{N}(\Gamma^\dagger \varepsilon'; \mathcal{F}, \rho_{\mathbf{X}}) \leq T \log |\mathcal{A}| + \max_{\alpha \in \mathcal{A}^T} \log \mathcal{N}(\Gamma^\dagger \varepsilon'; \mathcal{F}^\alpha, \rho_{\mathbf{X}}) \quad (17)$$

$$\leq T \log |\mathcal{A}| + \max_{\alpha \in \mathcal{A}^T} \left( \sum_{t=1}^T \log \mathcal{N}(\varepsilon'; \mathcal{H}, \rho_{\mathcal{Z}_{\mathcal{H}}^t}) + \sum_{\ell=1}^L \sum_{k=1}^{K_\ell} \log \mathcal{N}(\varepsilon'; \Psi_\ell, \rho_{\mathcal{Z}_\ell^k}) \right). \quad (18)$$

Till now, we have decomposed the covering number of  $\mathcal{F}^\alpha$  into product of covering numbers of all basic function sets  $\mathcal{H}, \Psi_\ell, \ell \in [L]$ . Next, following (Tripuraneni, Jordan, and Jin 2020), and the Sudakov minoration theorem (Theorem 5.30) and Lemma 5.5 in (Wainwright 2019), and recalling Definition 1, we have that for any  $\varepsilon' > 0$ ,

$$\begin{aligned}
\max_{\alpha \in \mathcal{A}^T} \sum_{t=1}^T \log \mathcal{N}(\varepsilon'; \mathcal{H}, \rho_{\mathcal{Z}_{\mathcal{H}}^t}) &\leq \max_{\alpha \in \mathcal{A}^T} \sum_{t=1}^T \left( \frac{2\sqrt{N}}{\varepsilon'} \widehat{\mathcal{G}}_{\mathcal{Z}_{\mathcal{H}}^t}(\mathcal{H}) \right)^2 \leq T \left( \frac{2\sqrt{N}}{\varepsilon'} \widetilde{\mathcal{G}}_N^{\mathcal{X}_{\mathcal{H}}}(\mathcal{H}) \right)^2, \\
\max_{\alpha \in \mathcal{A}^T} \sum_{k=1}^{K_\ell} \log \mathcal{N}(\varepsilon'; \Psi_\ell, \rho_{\mathcal{Z}_\ell^k}) &\leq \max_{\alpha \in \mathcal{A}^T} \sum_{k=1}^{K_\ell} \left( \frac{2\sqrt{|\mathcal{Z}_\ell^k|}}{\varepsilon'} \widehat{\mathcal{G}}_{\mathcal{Z}_\ell^k}(\Psi_\ell) \right)^2 \leq \max_{\alpha \in \mathcal{A}^T} \sum_{k=1}^{K_\ell} \left( \frac{2\sqrt{|\mathcal{Z}_\ell^k|}}{\varepsilon'} \widetilde{\mathcal{G}}_{|\mathcal{Z}_\ell^k|}^{\mathcal{X}_{\Psi_\ell}}(\Psi_\ell) \right)^2 \\
&\leq K_\ell \left( \frac{2\sqrt{2NT}}{\varepsilon'} \widetilde{\mathcal{G}}_{NT}^{\mathcal{X}_{\Psi_\ell}}(\Psi_\ell) \right)^2,
\end{aligned}$$

where the input spaces for  $\mathcal{H}$  and  $\Psi_\ell$  are  $\mathcal{X}_{\mathcal{H}} = \Psi_L \circ \dots \circ \Psi_1 \circ \mathcal{X}$ ,  $\mathcal{X}_{\Psi_\ell} = \Psi_{\ell-1} \circ \dots \circ \Psi_1 \circ \mathcal{X}$  for  $\ell > 1$  and  $\mathcal{X}_{\Psi_1} = \mathcal{X}$ . The last inequality is drawn from Lemma 3, which shows  $\sqrt{|\mathcal{Z}_\ell^k|} \widehat{\mathcal{G}}_{|\mathcal{Z}_\ell^k|}^{\mathcal{X}_{\Psi_\ell}}(\Psi_\ell) \leq \sqrt{2NT} \widetilde{\mathcal{G}}_{NT}^{\mathcal{X}_{\Psi_\ell}}(\Psi_\ell)$ . Since Definition 1 eliminates the input( $\mathbf{X}$ )-dependency, the inequalities hold for any valid inputs  $\mathbf{X}$ . In what follows, we drop the superscripts from the worst-case Gaussian complexities for cleaner exposition as they are clear from context. Then, setting  $\varepsilon' = \frac{u}{\Gamma^\dagger}$ , applying triangle inequality, we can obtain that for any  $\mathbf{X}$ ,

$$\sqrt{\log \mathcal{N}(u; \mathcal{F}, \rho_{\mathbf{X}})} \leq \sqrt{T \log |\mathcal{A}|} + \frac{2\Gamma^\dagger \sqrt{NT}}{u} \widetilde{\mathcal{G}}_N(\mathcal{H}) + \sum_{\ell=1}^L \frac{2\Gamma^\dagger \sqrt{2K_\ell NT}}{u} \widetilde{\mathcal{G}}_{NT}(\Psi_\ell). \quad (19)$$

Now it is time to combine everything together! Recall (14), (15) and (19). Since,  $D_{\mathbf{X}} \leq 2D_{\mathcal{X}}$  for any inputs  $\mathbf{X}$ , choosing  $\varepsilon = \frac{8D_{\mathcal{X}}}{NT}$ , we can obtain that with probability at least  $1 - \delta$ ,

$$\begin{aligned} \mathcal{R}_{\text{M}^2\text{TL}}(\hat{\mathbf{f}}) &\leq 6\Gamma\mathcal{G}_{NT}(\mathcal{F}) + 2\sqrt{\frac{\log \frac{2}{\delta}}{NT}} \\ &\leq 12\Gamma \left( \varepsilon + 32D_{\mathcal{X}}\sqrt{\frac{\log |\mathcal{A}|}{N}} + 32\Gamma^\dagger \left( \tilde{\mathcal{G}}_N(\mathcal{H}) + \sum_{\ell=1}^L \sqrt{2K_\ell}\tilde{\mathcal{G}}_{NT}(\Psi_\ell) \right) \int_{\varepsilon/4}^{2D_{\mathcal{X}}} \frac{1}{u} du \right) + 2\sqrt{\frac{\log \frac{2}{\delta}}{NT}} \\ &\leq 768\Gamma \left( \frac{D_{\mathcal{X}}}{NT} + D_{\mathcal{X}}\sqrt{\frac{\log |\mathcal{A}|}{N}} + \Gamma^\dagger \log NT \left( \tilde{\mathcal{G}}_N(\mathcal{H}) + \sum_{\ell=1}^L \sqrt{K_\ell}\tilde{\mathcal{G}}_{NT}(\Psi_\ell) \right) \right) + 2\sqrt{\frac{\log \frac{2}{\delta}}{NT}}. \end{aligned}$$

Till now, we have obtained the result for general  $\hat{\mathbf{f}}$ . Finally, consider the case that  $\hat{\mathbf{f}}$  might not utilize all the modules in the supernet. Let  $\hat{K}_\ell \leq K_\ell$  be the number of modules used by the empirical solution  $\hat{\mathbf{f}}$ . Applying Lemma 4, we can now replace  $\Phi$  with  $\Phi_{\text{used}}$  which replaces  $K_\ell$  with  $\hat{K}_\ell$  for  $\ell \in [L]$ , which concludes our final result.  $\blacksquare$

• **Developing an input-dependent bound.** In Theorem 1, we present the bound of Multipath MTL problem based on the worst-case Gaussian complexity. However, as shown in Definition 1, it computes the complexity of a function set by searching for the worst-case latent inputs, which ignores the data distribution and how the data collected as tasks. In the following argument, we present an input-based guarantee that bounds the excess risk of Multipath MTL problem tightly. To begin with, recall that  $\mathbf{X} = \{\mathbf{X}_t\}_{t=1}^T$  and  $\mathbf{X}_t = \{\mathbf{x}_{ti}\}_{i=1}^N$  denote the actual raw feature sets. Given inputs in  $T$  tasks, we can define the *empirical* worst-case Gaussian complexities of  $\mathcal{H}$  and  $\Psi_\ell, \ell \in [L]$  as follows.

$$\begin{aligned} C_{\mathbf{X}}^{\mathcal{H}} &= \max_{t \in [T]} \sup_{\mathbf{Z} \in \mathcal{Z}_t} \hat{\mathcal{G}}_{\mathbf{Z}}(\mathcal{H}), \quad \text{where } \mathcal{Z}_t = \Psi_L \circ \dots \circ \Psi_1(\mathbf{X}_t), \\ C_{\mathbf{X}}^{\Psi_\ell} &= \max_{\mathcal{I} \subset [T]} \sup_{\mathbf{Z} \in \mathcal{Z}_{\mathcal{I}}} \sqrt{\frac{|\mathcal{I}|}{T}} \hat{\mathcal{G}}_{\mathbf{Z}}(\Psi_\ell), \quad \text{where } \mathcal{Z}_{\mathcal{I}} = \bigcup_{t \in \mathcal{I}} \Psi_{\ell-1} \circ \dots \circ \Psi_1(\mathbf{X}_t), \end{aligned}$$

where  $\hat{\mathcal{G}}_{\mathbf{Z}}(\mathcal{H})$  and  $\hat{\mathcal{G}}_{\mathbf{Z}}(\Psi_\ell)$  are empirical Gaussian complexities and input spaces of  $\mathcal{H}$  and  $\Psi_\ell$  are corresponding to the raw input  $\mathbf{X}$ . Then, such statement provide another method to bound (18). That is, we have for any  $\varepsilon' > 0$ ,

$$\begin{aligned} \max_{\alpha \in \mathcal{A}^T} \sum_{t=1}^T \log \mathcal{N}(\varepsilon'; \mathcal{H}, \rho_{\mathbf{Z}_t^t}) &\leq \sum_{t=1}^T \left( \frac{2\sqrt{N}}{\varepsilon'} \max_{\alpha \in \mathcal{A}^T} \hat{\mathcal{G}}_{\mathbf{Z}_t^t}(\mathcal{H}) \right)^2 \leq T \left( \frac{2\sqrt{N}}{\varepsilon'} C_{\mathbf{X}}^{\mathcal{H}} \right)^2, \\ \max_{\alpha \in \mathcal{A}^T} \sum_{k=1}^{K_\ell} \log \mathcal{N}(\varepsilon'; \Psi_\ell, \rho_{\mathbf{Z}_\ell^k}) &\leq \sum_{k=1}^{K_\ell} \left( \frac{2\sqrt{NT}}{\varepsilon'} \max_{\alpha \in \mathcal{A}^T} \sqrt{\frac{|\mathbf{Z}_\ell^k|}{NT}} \hat{\mathcal{G}}_{\mathbf{Z}_\ell^k}(\Psi_\ell) \right)^2 \leq K_\ell \left( \frac{2\sqrt{NT}}{\varepsilon'} C_{\mathbf{X}}^{\Psi_\ell} \right)^2. \end{aligned}$$

The statements provided to prove Theorem 1 utilize the worst-case Gaussian complexity, and it bounds both empirical and population Gaussian complexities. Here,  $C_{\mathbf{X}}^{\mathcal{H}}$  and  $C_{\mathbf{X}}^{\Psi_\ell}$  depend on the input  $\mathbf{X}$ , and by construction, they are larger than their corresponding empirical complexities, however there is no guarantee that they will be larger than the corresponding population Gaussian complexities. Combining the result with (14), we can obtain that with probability at least  $1 - \delta$ ,

$$\mathcal{R}_{\text{M}^2\text{TL}}(\hat{\mathbf{f}}) \leq 384\Gamma \left( \frac{D_{\mathbf{X}}}{NT} + D_{\mathbf{X}}\sqrt{\frac{\log |\mathcal{A}|}{N}} + \Gamma^\dagger \log NT \left( C_{\mathbf{X}}^{\mathcal{H}} + \sum_{\ell=1}^L \sqrt{K_\ell} C_{\mathbf{X}}^{\Psi_\ell} \right) \right) + 6\sqrt{\frac{\log \frac{4}{\delta}}{NT}}, \quad (20)$$

where  $D_{\mathbf{X}} = \sup_{\mathbf{f}, \mathbf{f}' \in \mathcal{F}} \rho_{\mathbf{X}}(\mathbf{f}, \mathbf{f}')$ . Here we consider complexity of each task-specific head separately and bound it using the task with the largest head complexity ( $C_{\mathbf{X}}^{\mathcal{H}}$ ). As for the complexity of each layer, in the general case (as shown in Theorem 1), all the modules in the same layer share the same input space  $\mathcal{X}_{\Psi_\ell}$  by assuming raw input space  $\mathcal{X}$ , and because of Lemma 3, the sample complexity of  $\ell$ -th layer is bounded by  $\mathcal{O}\left(\sqrt{K_\ell}\tilde{\mathcal{G}}_{NT}(\Psi_\ell)\right)$ . When given actual training data  $\mathbf{X}$ , we need to search to find the worst-case cluster method of  $\ell$ -th layer, which results in  $C_{\mathbf{X}}^{\Psi_\ell}$ .

Below, we extend our theoretical result of Multipath MTL to two specific settings, vanilla MTL and hierarchical MTL.

**Corollary 2 (Vanilla MTL)** *Given the same data setting described in Section 2, consider a vanilla MTL problem as depicted in Figure 4a, which can be formulated as follows.*

$$\{\hat{h}_t\}_{t=1}^T, \hat{\phi} = \arg \min_{h_t \in \mathcal{H}, \phi \in \Phi} \frac{1}{NT} \sum_{t=1}^T \sum_{i=1}^N \ell(h_t \circ \phi(\mathbf{x}_{ti}), y_{ti}).$$

Suppose  $\mathcal{H}, \Phi$  are sets of  $\Gamma$ -Lipschitz functions with respect to Euclidean norm, and  $\ell(\cdot, y) : \mathbb{R} \times \mathbb{R} \rightarrow [0, 1]$  is also  $\Gamma$ -Lipschitz with respect to Euclidean norm. Define  $D_{\mathcal{X}} = \sup_{\mathbf{x} \in \mathcal{X}, h \in \mathcal{H}, \phi \in \Phi} |h \circ \phi(\mathbf{x})| < \infty$ . Let  $\mathcal{L}(\{h_t\}_{t=1}^T, \phi) = \mathbb{E}_{\mathcal{D}}[\ell(h_t \circ \phi(\mathbf{x}), y)]$  and  $\mathcal{L}^* = \min_{h_t \in \mathcal{H}, \phi \in \Phi} \mathbb{E}_{\mathcal{D}}[\ell(h_t \circ \phi(\mathbf{x}), y)]$ . Then we have that with probability at least  $1 - \delta$ ,

$$\mathcal{L}(\{\hat{h}_t\}_{t=1}^T, \hat{\phi}) - \mathcal{L}^* \leq 384\Gamma \left( \frac{D_{\mathcal{X}}}{NT} + (\Gamma + 1) \log NT \left( \tilde{\mathcal{G}}_N(\mathcal{H}) + \mathcal{G}_{NT}(\Phi) \right) \right) + 2\sqrt{\frac{\log \frac{2}{\delta}}{NT}}.$$

Here, the input space for  $\mathcal{H}$  is  $\Psi \times \mathcal{X}$ .

This corollary is consistent with (Tripuraneni, Jordan, and Jin 2020), and it can be simply deduced following the statement of Theorem 5, by setting  $L = 1, K_1 = 1$ . Since there is only one pathway selection,  $|\mathcal{A}| = 1$  and  $\log |\mathcal{A}| = 0$ . Here, the input space for representation  $\Phi$  is  $\mathcal{X}$ , and its complexity is shown in Gaussian complexity fashion.

**Corollary 3 (Hierarchical MTL)** Consider the hierarchical MTL problem depicted in Fig. 4c and consider a hierarchical supernet with degree  $K$ . Follow the same settings in Section 2. Suppose Assumptions 1&2 hold. Let  $\hat{\mathbf{f}}$  be the empirical solution of (M<sup>2</sup>TL). Let  $D_{\mathcal{X}} = \sup_{\mathbf{x} \in \mathcal{X}, h \in \mathcal{H}, \phi \in \Phi, \alpha \in \mathcal{A}} |h \circ \phi_{\alpha}(\mathbf{x})| < \infty$  and  $\Gamma^{\dagger} = \sum_{\ell=0}^L \Gamma^{\ell}$ . Then, with probability at least  $1 - \delta$ , the excess test risk in (2) obeys

$$\mathcal{R}_{M^2TL}(\hat{\mathbf{f}}) \leq 768\Gamma \left( \frac{D_{\mathcal{X}}}{NT} + D_{\mathcal{X}} \sqrt{\frac{(L-1) \log K}{N}} + \Gamma^{\dagger} \log NT \left( \tilde{\mathcal{G}}_N(\mathcal{H}) + \sum_{\ell=1}^L K^{\frac{\ell-1}{2}} \tilde{\mathcal{G}}_{NT}(\Psi_{\ell}) \right) \right) + 2\sqrt{\frac{\log \frac{2}{\delta}}{NT}}.$$

Here, the input spaces for  $\mathcal{H}$  and  $\Psi_{\ell}$  are  $\mathcal{X}_{\mathcal{H}} = \Psi_L \circ \dots \circ \Psi_1 \circ \mathcal{X}$ ,  $\mathcal{X}_{\Psi_{\ell}} = \Psi_{\ell-1} \circ \dots \circ \Psi_1 \circ \mathcal{X}$  for  $\ell > 1$ , and  $\mathcal{X}_{\Psi_1} = \mathcal{X}$ . Now if we consider a two-layer hierarchical representations as depicted in Fig. 1b, we can immediately obtain the result by setting  $L = 2$  ( $\Gamma^{\dagger} = 1 + \Gamma + \Gamma^2$ ). Then with probability at least  $1 - \delta$ ,

$$\mathcal{R}_{M^2TL}(\hat{\mathbf{f}}) \leq 768\Gamma \left( \frac{D_{\mathcal{X}}}{NT} + D_{\mathcal{X}} \sqrt{\frac{\log K}{N}} + \Gamma^{\dagger} \log NT \left( \tilde{\mathcal{G}}_N(\mathcal{H}) + \mathcal{G}_{NT}(\Psi_1) + \sqrt{K} \tilde{\mathcal{G}}_{NT}(\Psi_2) \right) \right) + 2\sqrt{\frac{\log \frac{2}{\delta}}{NT}}.$$

The result is consistent with Section 5, and proof can be immediately done by setting  $|\mathcal{A}| = K^{L-1}$  and  $K_{\ell} = K^{\ell-1}$  in Theorem 5. Here we observe that if the complexity of  $\Psi_{\ell}$  decreasing exponentially as  $\text{comp}(\Psi_{\ell}) \propto K^{-\frac{\ell}{2}}$ , then each layer has a constant complexity. We believe this and similar bounds can potentially provide guidelines on how we should design hierarchical supernet.

### B.3 Proof of Lemma 1

**Lemma 5 (Lemma 1 restated)** Recall  $\hat{\mathbf{f}}$  is the solution of (M<sup>2</sup>TL) and  $\hat{f}_t = \hat{h}_t \circ \hat{\phi}_{\hat{\alpha}_t}$  is the associated task- $t$  hypothesis. Define the excess risk of task  $t$  as  $\mathcal{R}_t(\hat{f}_t) = \mathcal{L}_t(\hat{f}_t) - \mathcal{L}_t^*$  where  $\mathcal{L}_t(f) = \mathbb{E}_{\mathcal{D}_t}[\hat{\mathcal{L}}_t(f)]$  is the population risk of task  $t$  and  $\mathcal{L}_t^*$  is the optimal achievable test risk for task  $t$  over  $\mathcal{F}$ . With probability at least  $1 - \delta - \mathbb{P}(\hat{\mathcal{L}}_{S_{\text{all}}}(\hat{\mathbf{f}}) \neq 0)$ , for all tasks  $t \in [T]$ ,

$$\mathcal{R}_t(\hat{f}_t) \lesssim \tilde{\mathcal{G}}_N(\mathcal{H}) + \sum_{\ell=1}^L \tilde{\mathcal{G}}_N(\Psi_{\ell}) + \sqrt{\frac{\log(2T/\delta)}{N}}. \quad (21)$$

**Proof** Let  $\mathcal{F}_{\text{IND}}$  be the hypothesis class of a single task induced by a pathway in the supernet. Since modules are same,  $\mathcal{F}_{\text{IND}}$  is same regardless of pathway. First, applying our main theorem (Thm 1) for a single supernet with  $K_{\ell} = 1$  (i.e. on  $\mathcal{F}_{\text{IND}}$ ), for a single task  $t$ , we end up with the uniform concentration guarantee, for all  $f \in \mathcal{F}_{\text{IND}}$ , with probability at least  $1 - \delta$ ,

$$|\hat{\mathcal{L}}_{S_t}(f) - \mathcal{L}_t(f)| \lesssim \tilde{\mathcal{G}}_N(\mathcal{H}) + \sum_{\ell=1}^L \tilde{\mathcal{G}}_N(\Psi_{\ell}) + \sqrt{\frac{\log(2/\delta)}{N}}.$$

Union bounding, for all  $f_t \in \mathcal{F}_{\text{IND}}, t \in [T]$ , with probability at least  $1 - \delta$ , we obtain

$$|\hat{\mathcal{L}}_{S_t}(f_t) - \mathcal{L}_t(f_t)| \lesssim \tilde{\mathcal{G}}_N(\mathcal{H}) + \sum_{\ell=1}^L \tilde{\mathcal{G}}_N(\Psi_{\ell}) + \sqrt{\frac{\log(2T/\delta)}{N}}. \quad (22)$$

Let us call this intersection event  $\mathcal{E}_{\text{all}}$ . Intersecting this with the events  $\min_{f_t \in \mathcal{F}_{\text{IND}}} \hat{\mathcal{L}}_{S_t}(f_t) = 0$  for  $t \in [T]$ , we exactly end up with (21). Thus, the statement is indeed what one would obtain by union bounding individualized training.

To proceed, we argue that same bound holds when solving (M<sup>2</sup>TL). We know (22) holds for all  $f_t$  chosen from  $\mathcal{F}_{\text{IND}}$ , therefore it holds for  $\hat{f}_t, t \in [T]$ . Consider its intersection with the event  $\mathbb{P}(\hat{\mathcal{L}}_{S_{\text{all}}}(\hat{\mathbf{f}}) \neq 0)$ . Given that  $\hat{\mathcal{L}}_{S_t}(\hat{f}_t) = 0$ , we obtain  $\mathcal{R}_t(\hat{f}_t) \leq \mathcal{L}_t(\hat{f}_t)$  upper bounded by the RHS of (22). ■

## B.4 Proof of Theorem 2

**Theorem 6 (Theorem 2 restated)** *Suppose Assumptions 1&2 hold. Let supernet  $\hat{\phi}$  be the solution of (M<sup>2</sup>TL) and  $\hat{f}_{\hat{\phi}}$  be the empirical minima of (TLOP) with respect to supernet  $\hat{\phi}$ . Let  $D_{\mathcal{X}} = \sup_{\mathbf{x} \in \mathcal{X}, \alpha \in \mathcal{A}, h \in \mathcal{H}_{\mathcal{T}}} |h \circ \hat{\phi}_{\alpha}(\mathbf{x})| < \infty$ . Then with probability at least  $1 - \delta$ ,*

$$\mathcal{R}_{\text{TLOP}}(\hat{f}_{\hat{\phi}}) \leq \text{Bias}_{\mathcal{T}}(\hat{\phi}) + 768\Gamma \left( \frac{D_{\mathcal{X}}}{M} + D_{\mathcal{X}} \sqrt{\frac{\log |\mathcal{A}|}{M}} + \log M \cdot \tilde{\mathcal{G}}_M(\mathcal{H}_{\mathcal{T}}) \right) + 2\sqrt{\frac{\log \frac{2}{\delta}}{M}},$$

where input space of  $\tilde{\mathcal{G}}_M(\mathcal{H}_{\mathcal{T}})$  is given by  $\{\hat{\phi}_{\alpha} \circ \mathcal{X} \mid \alpha \in \mathcal{A}\}$ .

**Proof** For short notation, let  $\mathcal{H} := \mathcal{H}_{\mathcal{T}}$ . We consider the transfer learning problem over a target task, with distribution  $\mathcal{D}_{\mathcal{T}}$  and training dataset  $\mathcal{S}_{\mathcal{T}} = \{(\mathbf{x}_i, y_i)\}_{i=1}^M$  with  $M$  samples i.i.d. drawn from  $\mathcal{D}_{\mathcal{T}}$ . Let  $\hat{\phi}$  and  $\phi^*$  denote the empirical and population solution of (M<sup>2</sup>TL). Then, we can recap the excess transfer learning risk

$$\mathcal{R}_{\text{TLOP}}(\hat{f}_{\hat{\phi}}) = \mathcal{L}_{\mathcal{T}}(\hat{f}_{\hat{\phi}}) - \mathcal{L}_{\mathcal{T}}^* = \underbrace{\mathcal{L}_{\mathcal{T}}(\hat{f}_{\hat{\phi}}) - \mathcal{L}_{\mathcal{T}}(f_{\hat{\phi}}^*)}_{\text{variance (a)}} + \underbrace{\mathcal{L}_{\mathcal{T}}(f_{\hat{\phi}}^*) - \mathcal{L}_{\mathcal{T}}^*}_{\text{supernet bias (b)}}.$$

Following Definition 2,  $b = \text{Bias}_{\mathcal{T}}(\hat{\phi})$ , and it remains to bound variance (a). Let  $\hat{f}_{\hat{\phi}} := (\hat{h}_{\hat{\phi}}, \hat{\alpha}_{\hat{\phi}})$  and  $f_{\hat{\phi}}^* := (h_{\hat{\phi}}^*, \alpha_{\hat{\phi}}^*)$ . For short notations, we remove the subscript  $\hat{\phi}$ , and we assume supernet  $\hat{\phi}$  is implied. Following the similar statements in Appendix B.2, we can decompose variance as follows.

$$a = \mathcal{L}_{\mathcal{T}}(\hat{f}) - \mathcal{L}_{\mathcal{T}}(f^*) = \underbrace{\mathcal{L}_{\mathcal{T}}(\hat{f}) - \hat{\mathcal{L}}_{\mathcal{T}}(\hat{f})}_c + \underbrace{\hat{\mathcal{L}}_{\mathcal{T}}(\hat{f}) - \hat{\mathcal{L}}_{\mathcal{T}}(f^*)}_d + \underbrace{\hat{\mathcal{L}}_{\mathcal{T}}(f^*) - \mathcal{L}_{\mathcal{T}}(f^*)}_e$$

where  $\mathcal{L}_{\mathcal{T}}(f) = \mathbb{E}_{\mathcal{D}_{\mathcal{T}}}[\ell(h \circ \hat{\phi}_{\alpha}(\mathbf{x}), y)]$  and  $\hat{\mathcal{L}}_{\mathcal{T}}(f) = \frac{1}{M} \sum_{i=1}^M \ell(h \circ \hat{\phi}_{\alpha}(\mathbf{x}_i), y_i)$  where  $f = (h, \alpha)$  and  $(\mathbf{x}_i, y_i) \in \mathcal{S}_{\mathcal{T}}$ . Since  $\hat{f}$  minimizes the training loss given  $\hat{\phi}$ ,  $d \leq 0$ . Let  $\mathbf{X}$  denote the input dataset, that is,  $\mathbf{X} = \{\mathbf{x}_i\}_{i=1}^M$ . Same as Inequality (14) in Appendix B.2, we derive the similar result that with probability at least  $1 - \delta$ ,

$$\begin{aligned} \mathcal{R}_{\text{TLOP}}(\hat{f}) &\leq \text{Bias}_{\mathcal{T}}(\hat{\phi}) + 6\Gamma \mathcal{G}_M(\mathcal{H} \circ \hat{\phi}(\mathcal{A})) + 2\sqrt{\frac{\log \frac{2}{\delta}}{M}}, \\ \text{and } \mathcal{R}_{\text{TLOP}}(\hat{f}) &\leq \text{Bias}_{\mathcal{T}}(\hat{\phi}) + 6\Gamma \hat{\mathcal{G}}_{\mathbf{X}}(\mathcal{H} \circ \hat{\phi}(\mathcal{A})) + 6\sqrt{\frac{\log \frac{2}{\delta}}{M}}, \end{aligned}$$

where  $\hat{\mathcal{G}}_{\mathbf{X}}(\mathcal{H} \circ \hat{\phi}(\mathcal{A})) = \mathbb{E}_{\mathbf{g}} \left[ \sup_{h \in \mathcal{H}, \alpha \in \mathcal{A}} \frac{1}{M} \sum_{i=1}^M g_i h \circ \hat{\phi}_{\alpha}(\mathbf{x}_i) \right]$  and  $\mathcal{G}_M(\mathcal{H} \circ \hat{\phi}(\mathcal{A})) = \mathbb{E}_{\mathcal{D}_{\mathcal{T}}} \left[ \hat{\mathcal{G}}_{\mathbf{X}}(\mathcal{H} \circ \hat{\phi}(\mathcal{A})) \right]$ . Following the Definition 4, let  $D = \sup_{h, h' \in \mathcal{H}, \alpha, \alpha' \in \mathcal{A}} \rho_{\mathbf{X}}(h \circ \hat{\phi}_{\alpha}, h' \circ \hat{\phi}_{\alpha'}) \leq 2D_{\mathcal{X}}$ . By applying the Dudley's theorem, and following the same statements in Appendix B.2, we obtain that given any  $\varepsilon \in [0, D]$

$$\hat{\mathcal{G}}_{\mathbf{X}}(\mathcal{H} \circ \hat{\phi}(\mathcal{A})) \leq 2\varepsilon + \frac{32}{\sqrt{M}} \int_{\varepsilon/4}^D \sqrt{\log \mathcal{N}(u; \mathcal{H} \circ \hat{\phi}(\mathcal{A}), \rho_{\mathbf{X}})} du.$$

Now we need to decompose the covering number of  $\mathcal{H} \circ \hat{\phi}(\mathcal{A})$  into the covering numbers of separate hypothesis sets  $\mathcal{H}$  and  $\mathcal{A}$ . For short notations, let  $\mathcal{H}(\mathcal{A}) := \mathcal{H} \circ \hat{\phi}(\mathcal{A})$  and  $\mathcal{H}(\alpha) := \mathcal{H} \circ \hat{\phi}_{\alpha}$ , and we omit the subscript  $\mathbf{X}$  from  $\rho$ . Since pathway set  $\mathcal{A}$  is discrete with cardinality  $|\mathcal{A}|$ , the covering number of  $\mathcal{H}(\mathcal{A})$  is the product of covering number of  $\mathcal{H}(\alpha)$  for all  $\alpha \in \mathcal{A}$ , and can be bounded by the  $|\mathcal{A}|$  times product of the worst-case covering number of  $\mathcal{H}(\alpha)$ , that is  $\mathcal{N}(u; \mathcal{H}(\mathcal{A}), \rho) = \prod_{\alpha \in \mathcal{A}} \mathcal{N}(u; \mathcal{H}(\alpha), \rho) \leq \max_{\alpha \in \mathcal{A}} \mathcal{N}^{|\mathcal{A}|}(u; \mathcal{H}(\alpha), \rho)$ . Logarithm of it results in  $\log \mathcal{N}(u; \mathcal{H}(\mathcal{A}), \rho) \leq \log |\mathcal{A}| + \max_{\alpha \in \mathcal{A}} \log \mathcal{N}(u; \mathcal{H}(\alpha), \rho)$ . Now let  $\mathbf{Z}_{\alpha} = \hat{\phi}_{\alpha}(\mathbf{X}) = \{\hat{\phi}_{\alpha}(\mathbf{x}_i) : \mathbf{x}_i \in \mathbf{X}\}$ , which is the set of latent inputs of prediction head. Then for any given  $\alpha \in \mathcal{A}$ ,

$$\rho_{\mathbf{X}}(h \circ \hat{\phi}_{\alpha}, h' \circ \hat{\phi}_{\alpha}) = \sqrt{\frac{1}{M} \sum_{i=1}^M (h \circ \hat{\phi}_{\alpha}(\mathbf{x}_i) - h' \circ \hat{\phi}_{\alpha}(\mathbf{x}_i))^2} = \sqrt{\frac{1}{M} \sum_{i=1}^M (h(\mathbf{z}_i) - h'(\mathbf{z}_i))^2} = \rho_{\mathbf{Z}_{\alpha}}(h, h'),$$

where  $\mathbf{z}_i = \hat{\phi}_{\alpha}(\mathbf{x}_i)$  and then  $\mathbf{Z}_{\alpha} = \{\mathbf{z}_1, \dots, \mathbf{z}_M\}$ . Such equality states that if pathway  $\alpha$  is fixed,  $u$ -cover of head  $\mathcal{H}$  results in  $u$ -cover of the prediction function, and simply,  $\mathcal{N}(u; \mathcal{H}(\alpha), \rho_{\mathbf{X}}) = \mathcal{N}(u; \mathcal{H}, \rho_{\mathbf{Z}_{\alpha}})$ . Next, following the same statements

in Appendix B.2, if we utilize the Sudakov minoration theorem in (Wainwright 2019), we obtain  $\sqrt{\log \mathcal{N}(u; \mathcal{H}, \rho_{\mathbf{Z}_\alpha})} \leq \frac{2\sqrt{M}}{u} \widehat{\mathcal{G}}_{\mathbf{Z}_\alpha}(\mathcal{H})$ . Finally, combining all we have together obtains

$$\begin{aligned} \widehat{\mathcal{G}}_{\mathbf{X}}(\mathcal{H}(\mathcal{A})) &\leq 2\varepsilon + \frac{32}{\sqrt{M}} \int_{\varepsilon/4}^D \sqrt{\log \mathcal{N}(u; \mathcal{H}(\mathcal{A}), \rho_{\mathbf{X}})} du \leq 2\varepsilon + 32D \sqrt{\frac{\log |\mathcal{A}|}{M}} + 64 \max_{\alpha \in \mathcal{A}} \widehat{\mathcal{G}}_{\mathbf{Z}_\alpha}(\mathcal{H}) \int_{\varepsilon/4}^D \frac{1}{u} du \\ &\leq 2\varepsilon + 32D \sqrt{\frac{\log |\mathcal{A}|}{M}} + 64 \log \frac{4D}{\varepsilon} \max_{\alpha \in \mathcal{A}} \widehat{\mathcal{G}}_{\mathbf{Z}_\alpha}(\mathcal{H}) \leq 64 \left( \frac{D}{M} + D \sqrt{\frac{\log |\mathcal{A}|}{M}} + \log M \max_{\alpha \in \mathcal{A}} \widehat{\mathcal{G}}_{\mathbf{Z}_\alpha}(\mathcal{H}) \right), \end{aligned}$$

by choosing  $\varepsilon = \frac{4D}{M}$ .

• **Input-dependent bound.** If we define the worst case *empirical* Gaussian complexity based on the raw input data  $\mathbf{X}$  and given supernet  $\hat{\phi}$ , that is  $C_{\mathbf{X}}^{\mathcal{H}} := \max_{\alpha \in \mathcal{A}} \widehat{\mathcal{G}}_{\mathbf{Z}_\alpha}(\mathcal{H})$ , where  $\mathbf{Z}_\alpha$  shows as above with respect to  $\hat{\phi}$  and  $\alpha$ , we have that with probability at least  $1 - \delta$ ,

$$\mathcal{R}_{\text{TLOP}}(\hat{f}_{\hat{\phi}}) \leq \text{Bias}_{\mathcal{T}}(\hat{\phi}) + 384\Gamma \left( \frac{D}{M} + D \sqrt{\frac{\log |\mathcal{A}|}{M}} + \log M \cdot C_{\mathbf{X}}^{\mathcal{H}} \right) + 6 \sqrt{\frac{\log \frac{4}{\delta}}{M}}.$$

Furthermore, let input space be  $\mathcal{X}$ . If we define the worst case Gaussian complexity independent to the specific training dataset and supernet, that is,  $\widetilde{\mathcal{G}}_M^{\mathcal{X}_\mathcal{H}}(\mathcal{H}) := \sup_{\mathbf{Z} \in \mathcal{X}_\mathcal{H}^M} \widehat{\mathcal{G}}_{\mathbf{X}}(\mathcal{H})$ , where  $\mathcal{X}_\mathcal{H} = \{\hat{\phi}_\alpha \circ \mathcal{X} | \alpha \in \mathcal{A}\}$ , then we have that

$$\mathcal{G}_M(\mathcal{H}(\mathcal{A})) \leq 64 \left( \frac{D_{\mathcal{X}}}{M} + D_{\mathcal{X}} \sqrt{\frac{\log |\mathcal{A}|}{M}} + \log M \cdot \widetilde{\mathcal{G}}_M^{\mathcal{X}_\mathcal{H}}(\mathcal{H}) \right),$$

which leads to the result that with probability at least  $1 - \delta$ ,

$$\mathcal{R}_{\text{TLOP}}(\hat{f}_{\hat{\phi}}) \leq \text{Bias}_{\mathcal{T}}(\hat{\phi}) + 384\Gamma \left( \frac{D_{\mathcal{X}}}{M} + D_{\mathcal{X}} \sqrt{\frac{\log |\mathcal{A}|}{M}} + \log M \cdot \widetilde{\mathcal{G}}_M(\mathcal{H}) \right) + 2 \sqrt{\frac{\log \frac{2}{\delta}}{M}}.$$

Here input space of  $\mathcal{H}$  is given by  $\mathcal{X}_\mathcal{H} = \{\hat{\phi}_\alpha \circ \mathcal{X} | \alpha \in \mathcal{A}\}$ . ■

## B.5 End-to-End Transfer Learning

In this section, we present an end-to-end transfer learning guarantee based on task diversity. We start with two useful definitions: supernet distance and task diversity. Here, supernet distance has been mentioned in Section 3.2 and following provides the intact definition. It measures the performance gap of two supernets. Similar to the previous work (Chen et al. 2021; Tripuraneni, Jordan, and Jin 2020; Xu and Tewari 2021), we define task diversity in Definition 6. It captures the similarity of target task to source tasks over a supernet by comparing their representation distance over it. Finally, using the task diversity argument, we can immediately obtain the theoretical guarantee for transfer learning risk.

**Definition 5 (Supernet Distance)** Consider a transfer learning with optimal pathway (TLOP) problem. Recall the definitions  $\mathcal{D}_{\mathcal{T}}$  and  $\mathcal{H}_{\mathcal{T}}$  stated in Section 2. Given two supernets  $\phi$  and  $\phi'$ , define the supernet/representation distance of  $\phi$  from  $\phi'$  for a target  $\mathcal{T}$  as

$$\text{Dist}_{\mathcal{T}}(\phi; \phi') = \text{Bias}_{\mathcal{T}}(\phi) - \text{Bias}_{\mathcal{T}}(\phi') = \min_{h \in \mathcal{H}_{\mathcal{T}}, \alpha \in \mathcal{A}} \mathcal{L}_{\mathcal{T}}(h \circ \phi_\alpha) - \min_{h \in \mathcal{H}_{\mathcal{T}}, \alpha \in \mathcal{A}} \mathcal{L}_{\mathcal{T}}(h \circ \phi'_\alpha).$$

Here, we do not restrict the supernet distance to target task  $\mathcal{T}$  only. Given source task  $t \in [T]$ , we can still define the corresponding supernet distance of  $\phi$  from  $\phi'$  as

$$\text{Dist}_t(\phi; \phi') = \min_{h \in \mathcal{H}, \alpha \in \mathcal{A}} \mathcal{L}_t(h \circ \phi_\alpha) - \min_{h \in \mathcal{H}, \alpha \in \mathcal{A}} \mathcal{L}_t(h \circ \phi'_\alpha), \quad (23)$$

and the hypothesis set for head is  $\mathcal{H}$  instead.

**Definition 6 (Task Diversity)** For any supernets  $\phi$  and  $\phi'$ , given  $T$  source tasks with distribution  $(\mathcal{D}_t)_{t=1}^T$  and a target task with distribution  $\mathcal{D}_{\mathcal{T}}$ , we say that the source tasks are  $(\nu, \epsilon)$ -diverse over the target task for a supernet  $\phi'$  if for any  $\phi \in \Phi$ ,

$$\text{Dist}_{\mathcal{T}}(\phi; \phi') \leq \left( \frac{1}{T} \sum_{t=1}^T \text{Dist}_t(\phi; \phi') \right) / \nu + \epsilon,$$

where we assume that head hypothesis sets  $\mathcal{H}$ ,  $\mathcal{H}_{\mathcal{T}}$  are implied for source and target distances.

**Theorem 7 (End-to-end transfer learning)** *Suppose Assumption 1&2 hold. Let supernet  $\hat{\phi}$  and  $\phi^*$  be the empirical and population solutions of (M<sup>2</sup>TL) and  $\hat{f}_{\hat{\phi}}$  be the empirical minima of (TLOP) with respect to supernet  $\hat{\phi}$ . Assume the source tasks used in Multipath MTL phase are  $(\nu, \epsilon)$ -diverse over target task  $\mathcal{T}$  for the optimal supernet  $\phi^*$ . Then with probability at least  $1 - 2\delta$ ,*

$$\mathcal{R}_{\text{TLOP}}(\hat{f}_{\hat{\phi}}) \lesssim \text{Bias}_{\mathcal{T}}(\phi^*) + \frac{1}{\nu} \left( \tilde{\mathcal{G}}_N(\mathcal{H}) + \sum_{\ell=1}^L \sqrt{\hat{K}_{\ell}} \tilde{\mathcal{G}}_{NT}(\Psi_{\ell}) + \sqrt{\frac{\log |\mathcal{A}|}{N}} \right) + \sqrt{\frac{\log |\mathcal{A}|}{M}} + \tilde{\mathcal{G}}_M(\mathcal{H}_{\mathcal{T}}) + \frac{1}{\nu} \sqrt{\frac{\log \frac{2}{\delta}}{NT}} + \sqrt{\frac{\log \frac{2}{\delta}}{M}} + \epsilon.$$

Here, the input spaces for  $\mathcal{H}$ ,  $\Psi_{\ell}$  and  $\mathcal{H}_{\mathcal{T}}$  are same to the statements in Theorem 1 and Theorem 2.

**Proof** Recall Theorem 6. To state end-to-end transfer learning risk, we need to bound supernet bias  $\text{Bias}_{\mathcal{T}}(\hat{\phi})$ . Following Definition 5, we have that  $\text{Bias}_{\mathcal{T}}(\hat{\phi}) = \text{Dist}_{\mathcal{T}}(\hat{\phi}; \phi^*) + \text{Bias}_{\mathcal{T}}(\phi^*)$ . Next, from Definition 6, since we assume source tasks are  $(\nu, \epsilon)$ -diverse over target task  $\mathcal{T}$  for the supernet  $\phi^*$ , we can obtain  $\text{Dist}_{\mathcal{T}}(\hat{\phi}; \phi^*) \leq \left( \frac{1}{T} \sum_{t=1}^T \text{Dist}_t(\hat{\phi}; \phi^*) \right) / \nu + \epsilon$ . To process, following (23), we have

$$\begin{aligned} \frac{1}{T} \sum_{t=1}^T \text{Dist}_t(\hat{\phi}; \phi^*) &= \frac{1}{T} \sum_{t=1}^T \left( \min_{h \in \mathcal{H}, \alpha \in \mathcal{A}} \mathcal{L}_t(h \circ \hat{\phi}_{\alpha}) - \min_{h \in \mathcal{H}, \alpha \in \mathcal{A}} \mathcal{L}_t(h \circ \phi_{\alpha}^*) \right) \\ &\leq \frac{1}{T} \sum_{t=1}^T \left( \mathcal{L}_t(\hat{h}_t \circ \hat{\phi}_{\hat{\alpha}_t}) - \mathcal{L}_t(h_t^* \circ \phi_{\alpha_t^*}^*) \right) = \mathcal{L}_{\bar{\mathcal{D}}}(\hat{f}) - \mathcal{L}_{\bar{\mathcal{D}}}^*(f) = \mathcal{R}_{\text{M}^2\text{TL}}(\hat{f}). \end{aligned}$$

Here,  $(\{\hat{h}_t, \hat{\alpha}_t\}_{t=1}^T, \hat{\phi})$  and  $(\{h_t^*, \alpha_t^*\}_{t=1}^T, \phi^*)$  are the empirical and population solutions of (M<sup>2</sup>TL), and we set  $\hat{f} := (\{\hat{h}_t, \hat{\alpha}_t\}_{t=1}^T, \hat{\phi})$ . The inequality term holds from the fact that: 1)  $\min_{h \in \mathcal{H}, \alpha \in \mathcal{A}} \mathcal{L}_t(h \circ \hat{\phi}_{\alpha}) \leq \mathcal{L}_t(\hat{h}_t \circ \hat{\phi}_{\hat{\alpha}_t})$ , and 2)  $\min_{h \in \mathcal{H}, \alpha \in \mathcal{A}} \mathcal{L}_t(h \circ \phi_{\alpha}^*) = \mathcal{L}_t(h_t^* \circ \phi_{\alpha_t^*}^*)$  since  $h_t^*$  and  $\alpha_t^*$  can be seen as the optimal solutions given supernet  $\phi^*$ . Combining them together with Theorem 1 and Theorem 2 completes the proof.  $\blacksquare$

## C Multipath MTL under Subexponential Loss Functions

The goal of this section is proving an MTL result under unbounded loss functions (e.g. least-squares). The high-level proof strategy is essentially a simplified version of proof of Theorem 1, where we use a more naive covering argument for parametric classes that have  $\mathcal{O}(\log(1/\epsilon))$  covering numbers. For this reason, we will make some simplifications in the proof to avoid repetitions. Instead, we will highlight key differences such as how the concentration argument changes due to unbounded losses. We first make the following assumptions.

**Assumption 7** *For any task distribution  $(\mathbf{x}, y) \sim \mathcal{D}_t$  and for any task hypothesis  $f_t \in \mathcal{F}_t$  (induced by  $\alpha_t, h_t, \phi$ ), we have that  $\ell(y, f_t(\mathbf{x}))$  is a  $\Xi$  subexponential random variable for some  $\Xi > 0$ . Additionally, assume loss is a  $\Gamma > 0$  Lipschitz function of  $\hat{y} = f_t(\mathbf{x})$ .*

We also assume a standard covering assumption. Note that, unlike the proof of Theorem 1, we focus on parametric classes and use data-agnostic covers.

**Assumption 8** *All modules  $\psi_{\ell}^k \in \Psi_{\ell}$  are  $\Gamma$  Lipschitz and map  $\psi_{\ell}^k : \mathbb{R}^{p_{\ell}-1} \rightarrow \mathbb{R}^{p_{\ell}}$ . For the sake of simplicity assume  $\psi_{\ell}^k(0) = 0$  (e.g. neural net layer with ReLU activation). Additionally, for any Euclidean ball of radius  $R$ , the covering dimension of  $\Psi_{\ell}$  follows the parametric classes, namely,*

$$\mathcal{N}(\epsilon; \Psi_{\ell}, R) \leq d_{\ell} \log\left(\frac{3R}{\epsilon}\right),$$

where  $d_{\ell}$  is the covering dimension of  $\Psi_{\ell}$ . Verbally, there exists a cover  $\Psi_{\ell}^{\epsilon}$ ,  $|\Psi_{\ell}^{\epsilon}| \leq \mathcal{N}(\epsilon; \Psi_{\ell}, R)$ , such that for any  $\|\mathbf{x}\| \leq R$  and for any  $\psi \in \Psi_{\ell}$ , there exists  $\psi' \in \Psi_{\ell}^{\epsilon}$  such that  $\|\psi'(\mathbf{x}) - \psi(\mathbf{x})\| \leq \epsilon$ . Additionally, let head  $\mathcal{H}$  be  $\Gamma$  Lipschitz and  $d_{\mathcal{H}}$  be the covering dimension for  $\mathcal{H}$ .

**Theorem 8** *Suppose  $\mathcal{X} \subset \mathcal{B}^p(R)$  and Assumptions 7 and 8 hold. Suppose we have a (M<sup>2</sup>TL) problem with  $N_{\text{tot}}$  samples in total where all training samples are independent, however, task sample sizes are arbitrary<sup>1</sup>. Note that, in the specific setting of Theorem 1, we have  $N_{\text{tot}} = NT$  with identical sample sizes. Assume that  $N_{\text{tot}} \gtrsim \text{DoF}(\mathcal{F}) \log(N_{\text{tot}}) + T \log |\mathcal{A}|$ . Declare population risk  $\mathcal{L}_{\bar{\mathcal{D}}}(\mathbf{f}) = \mathbb{E}[\hat{\mathcal{L}}_{S_{\text{all}}}(\mathbf{f})]$ . We have that with probability at least  $1 - \delta$ , for all Multipath hypothesis  $\mathbf{f} \in \mathcal{F}$*

$$|\hat{\mathcal{L}}_{S_{\text{all}}}(\mathbf{f}) - \mathcal{L}_{\bar{\mathcal{D}}}(\mathbf{f})| \lesssim \Xi \sqrt{\frac{L \cdot \text{DoF}(\mathcal{F}) + T \log |\mathcal{A}| + \log(2/\delta)}{N_{\text{tot}}}}. \quad (24)$$

<sup>1</sup>In words, tasks don't have to have  $N$  samples each. We will simply control the gap between empirical and population. If tasks have different sizes, then their population weights will similarly change.

The right hand side bounds similarly hold for the excess risk  $\mathcal{R}_{M^2TL}(\hat{\mathbf{f}})$  where  $\hat{\mathbf{f}}$  is the ERM solution.  $\lesssim$  subsumes the logarithmic dependence on  $R, \Gamma, L, N_{\text{tot}}$ . The exact bound is below (39). Finally, if we solve (M<sup>2</sup>TL) with fixed pathway choices (rather than searching over  $\mathcal{A}$ ), with same probability and assuming  $N_{\text{tot}} \gtrsim \text{DoF}(\mathcal{F}) \log(N_{\text{tot}})$  we have the simplified bound

$$|\hat{\mathcal{L}}_{S_{\text{all}}}(\mathbf{f}) - \mathcal{L}_{\mathcal{D}}(\mathbf{f})| \lesssim \Xi \sqrt{\frac{L \cdot \text{DoF}(\mathcal{F}) + \log(2/\delta)}{N_{\text{tot}}}}. \quad (25)$$

The theorem is automatically applicable to loss functions bounded by  $\Xi > 0$ . Also note that, this theorem avoids exponential depth dependence compared to Theorem 1. This is primarily because of the strong coverability of parametric classes which (essentially) applies a log operation to the Lipschitz constant of  $\mathcal{F}$ .

**Proof** Let  $R_0 = R$  and note that, at the  $\ell$ th layer, the input(output) space has radius  $R_{\ell-1}(R_\ell)$  where  $R_\ell = \Gamma^\ell R$ . Let  $\Psi_\ell$  denote the hypothesis set of the modules of  $\ell$ th layer. Fix an  $\varepsilon$  cover  $\mathcal{F}_\varepsilon$  for the sets  $(\Psi_\ell^{K_\ell})_{\ell=1}^L, \mathcal{H}^T$  and  $\mathcal{A}^T$ , such that  $\Psi_\ell$  is covered according to its input space radius  $R_{\ell-1}$  with resolution  $\varepsilon_\ell = \frac{\varepsilon}{\Gamma^{L-\ell+1}}$ , where  $\ell$  is layer depth and prediction head is layer  $L+1$ . This implies that

$$\log |\mathcal{F}_\varepsilon| \leq T d_{\mathcal{H}} \log \frac{3R_L}{\varepsilon} + T \log |\mathcal{A}| + \sum_{\ell=1}^L K_\ell d_\ell \log \frac{3R_{\ell-1}}{\varepsilon_\ell} \quad (26)$$

$$= (T d_{\mathcal{H}} + \sum_{\ell=1}^L K_\ell d_\ell) \log \frac{3R_L}{\varepsilon} + T \log |\mathcal{A}| \quad (27)$$

$$\leq \text{DoF}(\mathcal{F}) \log \frac{3R_L}{\varepsilon} + T \log |\mathcal{A}| \quad (28)$$

$$= \text{DoF}(\mathcal{F}) \left( \log \frac{3R}{\varepsilon} + L \log \Gamma \right) + T \log |\mathcal{A}|. \quad (29)$$

• **Step 1: Union bound over the cover.** We now show a uniform concentration argument over this cover. Since each sample is independent of others and each loss is  $\Xi$  subexponential, using subexponential Bernstein inequality (e.g. Prop 5.16 of (Vershynin 2010)), we have that

$$\mathbb{P} \left( |\hat{\mathcal{L}}_{S_{\text{all}}}(\mathbf{f}) - \mathcal{L}_{\mathcal{D}}(\mathbf{f})| \geq \frac{t}{\sqrt{N_{\text{tot}}}} \right) \leq 2 \exp \left( -c \min \left\{ \frac{t^2}{\Xi^2}, \frac{t\sqrt{N_{\text{tot}}}}{\Xi} \right\} \right)$$

Let  $\varepsilon = \frac{1}{\Gamma^{L+1} N_{\text{tot}}}$ , and recall that we assumed  $N_{\text{tot}} \gtrsim \log |\mathcal{F}_\varepsilon| + \tau$ . Now, setting  $t \propto \sqrt{\log |\mathcal{F}_\varepsilon| + \tau}$  and union bounding over all  $\mathbf{f} \in \mathcal{F}_\varepsilon$ , we find that, uniformly over  $\mathcal{F}_\varepsilon$ ,

$$\mathbb{P} \left( |\hat{\mathcal{L}}_{S_{\text{all}}}(\mathbf{f}) - \mathcal{L}_{\mathcal{D}}(\mathbf{f})| \geq \Xi \sqrt{\frac{\log |\mathcal{F}_\varepsilon| + \tau}{N_{\text{tot}}}} \right) \leq 2e^{-\tau}. \quad (30)$$

• **Step 2: Perturbation analysis.** Now that covering analysis is done, we proceed with controlling the perturbation. Let  $\mathbf{f} \in \mathcal{F}$  be a Multipath MTL hypothesis. We choose  $\mathbf{f}' \in \mathcal{F}_\varepsilon$  such that:

- $\mathbf{f}'$  chooses the same pathways.
- $\mathbf{f}'$  chooses heads  $(h'_t)_{t=1}^T$  and modules  $((\psi_\ell^{K_\ell})_{k=1}^L)_{\ell=1}^L$  such that these hypotheses are  $\varepsilon$  close over their respective input spaces to the hypotheses of  $\mathbf{f}$  denoted by  $(h_t)_{t=1}^T$  and modules  $((\psi_\ell^k)_{k=1}^L)_{\ell=1}^L$ .

Fix an arbitrary  $\mathbf{x} \in \mathcal{X}$  and task  $t \in [T]$ . Set the short-hand notation  $\tilde{\psi}_\ell = \psi_\ell^{\alpha_t}$  and  $\tilde{\psi}'_\ell = \psi'^{\alpha_t}$ . Along the pathway  $\alpha_t$ , define the functions

$$f_t^\ell(\mathbf{x}) = \begin{cases} f_t(\mathbf{x}) & \text{if } \ell = L+1, \\ f'_t(\mathbf{x}) & \text{if } \ell = 0, \\ h'_t \circ \tilde{\psi}'_L \circ \dots \circ \tilde{\psi}'_{\ell+1} \circ \tilde{\psi}_\ell \circ \dots \circ \tilde{\psi}_1(\mathbf{x}) & \text{if } 1 \leq \ell \leq L. \end{cases}$$

Let  $\mathbf{x}_\ell = \tilde{\psi}_\ell \circ \dots \circ \tilde{\psi}_1(\mathbf{x})$ . Recall that,  $\Psi_\ell$  is covered with resolution  $\varepsilon_\ell$ . Now, through a standard perturbation decomposition,

we find that

$$|f_t(\mathbf{x}) - f'_t(\mathbf{x})| \leq \sum_{\ell=0}^L |f_t^{\ell+1}(\mathbf{x}) - f_t^\ell(\mathbf{x})| \quad (31)$$

$$\leq \sum_{\ell=0}^L |h'_t \circ \tilde{\psi}'_L \circ \dots \circ \tilde{\psi}'_{\ell+1}(\mathbf{x}_\ell) - h'_t \circ \tilde{\psi}'_L \circ \dots \circ \tilde{\psi}'_{\ell+1}(\mathbf{x}_\ell)| \quad (32)$$

$$\leq \sum_{\ell=0}^L \Gamma^{L-\ell} \varepsilon_{\ell+1} \quad (33)$$

$$= (L+1)\varepsilon. \quad (34)$$

This establishes that if tasks choose the same pathways, proposed  $\varepsilon$  cover ensures that for all  $\mathbf{x} \in \mathcal{X}$  and task  $t$ ,  $\mathcal{F}_\varepsilon$  is an  $(L+1)\varepsilon$  cover of  $\mathcal{F}$ . To conclude, using  $\Gamma$  Lipschitzness of the loss function, we obtain

$$\mathcal{L}_{\bar{\mathcal{D}}}(\mathbf{f}) - \mathcal{L}_{\bar{\mathcal{D}}}(\mathbf{f}') \leq \sup_{\mathbf{x}, t} |\ell(y, f_t(\mathbf{x})) - \ell(y, f'_t(\mathbf{x}))| \leq \sup_{\mathbf{x}, t} \Gamma |f_t(\mathbf{x}) - f'_t(\mathbf{x})| \leq \Gamma(L+1)\varepsilon \quad (35)$$

$$\widehat{\mathcal{L}}_{\mathcal{S}_{\text{all}}}(\mathbf{f}) - \widehat{\mathcal{L}}_{\mathcal{S}_{\text{all}}}(\mathbf{f}') \leq \sup_{\mathbf{x}, t} |\ell(y, f_t(\mathbf{x})) - \ell(y, f'_t(\mathbf{x}))| \leq \sup_{\mathbf{x}, t} \Gamma |f_t(\mathbf{x}) - f'_t(\mathbf{x})| \leq \Gamma(L+1)\varepsilon. \quad (36)$$

Combining with uniform concentration, we found that, for all  $\mathbf{f} \in \mathcal{F}$ , with probability  $1 - \delta$ ,

$$|\widehat{\mathcal{L}}_{\mathcal{S}_{\text{all}}}(\mathbf{f}) - \mathcal{L}_{\bar{\mathcal{D}}}(\mathbf{f})| \lesssim \Xi \sqrt{\frac{\log |\mathcal{F}_\varepsilon| + \log(2/\delta)}{N_{\text{tot}}}} + \Gamma(L+1)\varepsilon \quad (37)$$

$$\lesssim \Xi \sqrt{\frac{\text{DoF}(\mathcal{F}) (\log \frac{3R}{\varepsilon} + L \log \Gamma) + T \log |\mathcal{A}| + \log(2/\delta)}{N_{\text{tot}}}} + \Gamma(L+1)\varepsilon. \quad (38)$$

Recall that  $\varepsilon = \frac{1}{\Gamma(L+1)N_{\text{tot}}}$ , then we obtain the advertised uniform concentration guarantee

$$|\widehat{\mathcal{L}}_{\mathcal{S}_{\text{all}}}(\mathbf{f}) - \mathcal{L}_{\bar{\mathcal{D}}}(\mathbf{f})| \lesssim \Xi \sqrt{\frac{\text{DoF}(\mathcal{F}) (\log(3R\Gamma(L+1)N_{\text{tot}}) + L \log \Gamma) + T \log |\mathcal{A}| + \log(2/\delta)}{N_{\text{tot}}}}. \quad (39)$$

We get the simplified statement (24) after ignoring the log factors. Finally, (25) arises by repeating above argument step-by-step while ignoring  $|\mathcal{A}|$  term in (26).  $\blacksquare$

## D Proofs in Section 4

### D.1 A direct corollary of Theorem 1 to linear representations

We start with a lemma that controls the worst-case Gaussian complexity of linear models. The proof is standard and stated for completeness.

**Lemma 6 (Linear models)** *Let  $\mathcal{B} \subset \mathbb{R}^{d \times p}$  be a set of matrices with operator norm bounded by a constant  $C > 0$  and let  $\mathcal{X} \subset \mathcal{B}^p(R)$  (subset of  $\ell_2$  ball of radius  $R$ ). Then*

$$\tilde{\mathcal{G}}_n^{\mathcal{X}}(\mathcal{B}) \leq CR \sqrt{\frac{dp}{n}}.$$

**Proof** Set  $\mathbf{X}_\varepsilon = \sum_{i=1}^n \mathbf{x}_i \mathbf{g}_i^\top = \mathbf{X}^\top \mathbf{G}$  where  $\mathbf{X} \in \mathbb{R}^{n \times p}$  is dataset and  $\mathbf{G} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, 1) \in \mathbb{R}^{n \times d}$ . Applying Cauchy-Schwarz, we write

$$\tilde{\mathcal{G}}_n^{\mathcal{X}}(\mathcal{B}) = \frac{1}{n} \sup_{\mathbf{X} \in \mathcal{X}^n} \mathbb{E} \left[ \sup_{\mathbf{B} \in \mathcal{B}} \sum_{i=1}^n \mathbf{g}_i^\top \mathbf{B} \mathbf{x}_i \right] = \frac{1}{n} \sup_{\mathbf{X} \in \mathcal{X}^n} \mathbb{E} \left[ \sup_{\mathbf{B} \in \mathcal{B}} \text{trace}(\mathbf{X}_\varepsilon \mathbf{B}) \right] \quad (40)$$

$$\leq C \frac{\sqrt{p}}{n} \sup_{\mathbf{X} \in \mathcal{X}^n} \mathbb{E} [\|\mathbf{X}_\varepsilon\|_F] \leq C \frac{\sqrt{p}}{n} \sup_{\mathbf{X} \in \mathcal{X}^n} \sqrt{\mathbb{E} [\|\mathbf{X}^\top \mathbf{G}\|_F^2]} \leq CR \sqrt{\frac{dp}{n}}. \quad (41)$$

$\blacksquare$



**Corollary 4** Suppose Assumptions 2&3 hold and input set  $\mathcal{X} \subset \mathcal{B}^p(c\sqrt{p})^2$  for a constant  $c > 0$ . Let  $\hat{\mathbf{f}}$  be empirical solution of (4). Then, with probability at least  $1 - \delta$ ,

$$\mathcal{R}_{M^2TL}(\hat{\mathbf{f}}) \lesssim \sqrt{\frac{p \cdot \text{DoF}(\mathcal{F})}{NT}} + \sqrt{\frac{\log |\mathcal{A}|}{N} + \frac{\log(2/\delta)}{NT}},$$

where  $\text{DoF}(\mathcal{F}) = T \cdot p_L + \sum_{\ell=1}^L K_\ell \cdot p_\ell \cdot p_{\ell-1}$  is the total number of trainable parameters in  $\mathcal{F}$ .

**Proof** This proof is immediately done by following Theorem 1 and Lemma 6. Since  $\mathcal{X} \subset \mathcal{B}^p(c\sqrt{p})$ , and we assume  $\Psi_\ell, \ell \in [L]$  have bounded operator norm  $C$ , for each layer, the input space  $\mathcal{X}_{\Psi_\ell} \subset \mathcal{B}^{p_{\ell-1}}(C^{\ell-1}c\sqrt{p})$  and then following Lemma 6,  $\tilde{\mathcal{G}}_{NT}(\Psi_\ell) \leq C^\ell c\sqrt{p} \sqrt{\frac{p_\ell p_{\ell-1}}{NT}}$ ,  $\ell \in [L]$ . Since  $\mathcal{H} = \mathcal{B}^{p_L}(C)$  and  $\mathcal{X}_{\mathcal{H}} \subset \mathcal{B}^{p_L}(C^L c\sqrt{p})$ , we have  $\tilde{\mathcal{G}}_N(\mathcal{H}) \leq C^{L+1} c\sqrt{p} \sqrt{\frac{p_L}{N}}$ . Then we obtain

$$\tilde{\mathcal{G}}_N(\mathcal{H}) + \sum_{\ell=1}^L \sqrt{K_\ell} \tilde{\mathcal{G}}_{NT}(\Psi_\ell) \lesssim c\sqrt{p} \cdot \sqrt{\frac{C^{L+1} \cdot T \cdot p_L + \sum_{\ell=1}^L C^\ell \cdot K_\ell \cdot p_\ell \cdot p_{\ell-1}}{NT}}.$$

Combining it with Theorem 1 finishes the proof.  $\blacksquare$

## D.2 Proof of Theorem 3

Corollary 4 directly follows by applying Theorem 1 to the linear representation setting, and therefore  $\lesssim$  subsumes dependencies on  $\log NT$  and  $\Gamma^L$ . Instead in Theorem 3 we establish a tighter bound for parametric hypothesis classes and the sample complexity is only logarithmic in the input space radius  $R$  ( $R = c\sqrt{p}$  in Corollary 4) and linearly dependent on the number of layers  $L$ .

**Proof** The theorem is a direct application of Theorem 8 after verifying the assumptions. First, bounded loss function  $\ell : \mathbb{R} \times \mathbb{R} \rightarrow [0, 1]$  implies it is sub-exponential, which verifies Assumption 7. Second, all module/head functions have bounded spectral/Euclidean norms, which verifies Assumption 8. One remark (compared to Theorem 8) is that, since the loss function is bounded, by applying Hoeffding's inequality, (30) holds without enforcing a lower bound constraint on  $N_{\text{tot}}$ .  $\blacksquare$

## D.3 The Need for Well-Populated Source Tasks in Assumption 5

**Lemma 7** Consider a weaker version of Assumption 5 where we enforce  $\Sigma_\alpha \succeq c\mathbf{I}_{p_L}$  over all clusters with two or more tasks<sup>3</sup> (i.e. only when  $\gamma_\alpha \geq 2/p_L$ ). Then, there exists a ((M<sup>2</sup>TL), (TLOP)) problem pair such that the excess transfer learning risk obeys  $\mathcal{R}_{TLOP}(\hat{\mathbf{f}}_\phi) \geq 1$  as soon as  $N \geq p$ .

**Proof** The idea is packing supernet with isolated MTL tasks that are uncorrelated with target while achieving zero MTL risk. We consider a simple supernet construction where  $T$  tasks will be processed in parallel and all layers have exactly  $K_\ell = T$  modules. Specifically, task  $t$  will use the pathway  $\alpha_t = [t, t, \dots, t]$  by selecting  $t$ th module from each layer. This way each task will use a unique pathway and supernet will be fully occupied. Set noise level  $\sigma = 0$ . Observe that as soon as  $N \geq p$ ,  $\theta_t^*$  minimizes both empirical and population risks. Consequently, for any  $\|\hat{\mathbf{h}}_t\| = 1$ ,  $\mathbf{B}_{\alpha_t} = \mathbf{h}_t(\theta_t^*)^\top$  is a valid (and minimum norm) minimizer of empirical and population risks. Here, we highlight the minimum norm aspect because this solution is what gradient descent would converge during MTL phase (while we acknowledge the existence of infinitely-many solutions) (Ji and Telgarsky 2018). To wrap up the proof, suppose transfer task is orthogonal to all source tasks and observe that, regardless of the transfer prediction head  $\hat{\mathbf{h}}_\mathcal{T}$  and pathway choice  $t$ , we have

$$\begin{aligned} \mathcal{R}_{TLOP}(\hat{\mathbf{f}}_\phi) &= \mathbb{E} \left[ (y - \hat{\mathbf{f}}_\phi(\mathbf{x}))^2 \right] = \mathbb{E} \left[ (\boldsymbol{\theta}_\mathcal{T}^\top \mathbf{x} - \hat{\mathbf{h}}_\mathcal{T}^\top \mathbf{h}_t(\boldsymbol{\theta}_t^*)^\top \mathbf{x})^2 \right] \\ &\geq \|\boldsymbol{\theta}_\mathcal{T} - (\hat{\mathbf{h}}_\mathcal{T}^\top \mathbf{h}_t) \boldsymbol{\theta}_t^*\|^2 \geq \|\boldsymbol{\theta}_\mathcal{T}\|^2 = 1. \end{aligned}$$

This concludes the proof. We note that, if  $\sigma \neq 0$  same argument would work as  $N \rightarrow \infty$ . Additionally, through same argument with  $\sigma = 0$ , it can be observed that, a more general lower bound on excess transfer risk is  $\min_{t \in [T]} \|\boldsymbol{\theta}_\mathcal{T} - \boldsymbol{\theta}_t^*\|^2/2$ .  $\blacksquare$

<sup>2</sup>Observe that, this input space is rich enough to capture a random vector with  $\mathcal{O}(1)$  subgaussian norm. For instance, a standard normal vector would fall into this set with exponentially high probability as soon as  $c > 1$ .

<sup>3</sup>The relaxation is not enforcing anything on pathways containing a single task.

## D.4 Proof of Theorem 4 and Supporting Results

We start with a useful lemma to show excess risk of linear least squares problem with dependent noise.

**Lemma 8 (Linear least squares risk with dependent noise)** *Let  $\mathcal{S} = (\mathbf{x}_i, y_i)_{i=1}^n \stackrel{i.i.d.}{\sim} \mathcal{D}$  where  $y = \boldsymbol{\theta}^\top \mathbf{x} + z$  where  $\mathbf{x}$  is  $\mathcal{O}(1)$  subgaussian vector with isotropic covariance and  $z$  is  $\mathcal{O}(\sigma)$  subgaussian noise. Here, we assume that  $\mathbf{x}$  &  $z$  can be dependent, however, orthogonal (i.e.  $\mathbb{E}[\mathbf{x}z] = 0$ ). Let  $\mathbf{X} = [\mathbf{x}_1 \cdots \mathbf{x}_n]^\top \in \mathbb{R}^{n \times p}$  and  $\mathbf{X}^\dagger$  be the Moore-Penrose pseudoinverse of  $\mathbf{X}$ . Let  $\wedge$  be the minimum symbol. For  $n \geq Cp$  for a sufficiently large constant  $C \geq 1$ , the excess least squares risk and population-empirical risk gap of  $\hat{\boldsymbol{\theta}} = \mathbf{X}^\dagger \mathbf{y}$  is given by*

$$\mathcal{L}_{\mathcal{D}}(\hat{\boldsymbol{\theta}}) - \sigma^2 \leq C\sigma^2 \frac{p+t}{n} \quad \text{with probability at least } 1 - e^{-cn} - 2e^{-\sqrt{tn}\wedge t} \quad (42)$$

$$\mathcal{L}_{\mathcal{D}}(\hat{\boldsymbol{\theta}}) - \hat{\mathcal{L}}_{\mathcal{S}}(\hat{\boldsymbol{\theta}}) \leq C\sigma^2 \left( \frac{p}{n} + \sqrt{\frac{t}{n}} \right) \quad \text{with probability at least } 1 - 2e^{-cn} - 4e^{-\sqrt{tn}\wedge t}. \quad (43)$$

**Proof** Let  $\mathbf{z} = [z_1 \dots z_n]^\top$  and  $\sigma_{\min}(\cdot), \sigma_{\max}(\cdot)$  return the smallest and biggest singular value of a matrix. We can write

$$\mathcal{L}_{\mathcal{D}}(\hat{\boldsymbol{\theta}}) - \mathbb{E}[z^2] = \|\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}\|^2 = \|(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{z}\|^2 \leq \frac{\|\mathbf{X}^\top \mathbf{z}\|^2}{\sigma_{\min}(\mathbf{X})^4}.$$

Following (Vershynin 2010), we have  $\sqrt{2n} \geq \sigma_{\max}(\mathbf{X}), \sigma_{\min}(\mathbf{X}) \geq \sqrt{n/2}$  each with probability at least  $1 - e^{-cn}$ . The crucial term of interest is  $\|\mathbf{X}^\top \mathbf{z}\|$ . To control this, observe that  $\mathbf{X}^\top \mathbf{z} = \sum_{i=1}^n z_i \mathbf{x}_i$ . Since  $z_i \mathbf{x}_i$  is  $\mathcal{O}(\sigma)$ -subexponential (multiplication of two subgaussians), the summand  $\mathbf{X}^\top \mathbf{z}$  has a mixed subgaussian/subexponential tail. Specifically, it obeys (Oymak 2018, Lemma D.7)

$$\mathbb{P}(\|\mathbf{X}^\top \mathbf{z}\|^2 \gtrsim \sigma^2(p+t)n) \leq 2e^{-\sqrt{tn}\wedge t}.$$

Combining both, with advertised probability we establish the first claim.

$$\frac{\|\mathbf{X}^\top \mathbf{z}\|^2}{\sigma_{\min}(\mathbf{X})^4} \lesssim \sigma^2 \frac{p+t}{n}.$$

For the second claim, observe that

$$\hat{\mathcal{L}}_{\mathcal{S}}(\hat{\boldsymbol{\theta}}) - \frac{1}{n}\|\mathbf{z}\|^2 = \frac{1}{n}\|\mathbf{y} - \hat{\mathbf{y}}\|^2 - \frac{1}{n}\|\mathbf{z}\|^2 = \frac{1}{n}[\|(\mathbf{I} - \mathbf{X}\mathbf{X}^\dagger)\mathbf{z}\|^2 - \|\mathbf{z}\|^2] \quad (44)$$

$$= \frac{1}{n}\|\mathbf{X}\mathbf{X}^\dagger \mathbf{z}\|^2 \leq \frac{\sigma_{\max}(\mathbf{X})^2}{n} \|(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{z}\|^2 \quad (45)$$

$$\leq 2\|(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{z}\|^2 \lesssim \sigma^2 \frac{p+t}{n}. \quad (46)$$

Here, the first and second inequalities of last line hold with respective probabilities at least  $1 - e^{-cn}$  and  $1 - e^{-cn} - 2e^{-\sqrt{tn}\wedge t}$ . Additionally, since  $z^2$  is  $\mathcal{O}(\sigma^2)$ -subexponential,  $|\frac{1}{n}\|\mathbf{z}\|^2 - \mathbb{E}[z^2]| \lesssim \sigma^2 \sqrt{t/n}$  with probability at least  $1 - 2e^{-\sqrt{tn}\wedge t}$ . Combining all provides the final equation bounding the gap between empirical and population risks.  $\blacksquare$

Then we present the following lemma that converts an MTL guarantee into a transfer learning guarantee on a single subspace.

**Lemma 9** *Let  $\mathbf{B} \in \mathbb{R}^{r \times p}$  be a matrix with orthonormal rows and fix  $\{\mathbf{h}_t\}_{t=1}^T \in \mathbb{R}^r$  with unit covariance and declare distributions  $(\mathbf{x}, y) \sim \mathcal{D}_t$  obeying  $y = \mathbf{h}_t^\top \mathbf{B}\mathbf{x} + z$  with  $\mathbb{E}[z^2] = \sigma^2$  and  $\mathbb{E}[\mathbf{x}\mathbf{x}^\top] = \mathbf{I}_p$ . Form  $\mathbf{H} = [\mathbf{h}_1 \dots \mathbf{h}_T]^\top$  and assume  $C\frac{1}{r}\mathbf{I}_r \succeq \frac{1}{T}\mathbf{H}^\top \mathbf{H} \succeq c\frac{1}{r}\mathbf{I}_r$ . Now, for some  $\varepsilon > 0$ , suppose that  $\hat{\mathbf{f}} = (\hat{\mathbf{B}}, \{\hat{\mathbf{h}}_t\}_{t=1}^T)$  with orthonormal  $\hat{\mathbf{B}}$  achieves small population risk in average that is*

$$\mathcal{L}_{\mathcal{D}}(\hat{\mathbf{f}}) - \mathcal{L}_{\mathcal{D}}(\mathbf{f}_*) = \frac{1}{T} \sum_{t=1}^T \mathbb{E}_{\mathcal{D}_t}[(\mathbf{h}_t^\top \mathbf{B}\mathbf{x} - \hat{\mathbf{h}}_t^\top \hat{\mathbf{B}}\mathbf{x})^2] \leq \varepsilon$$

where  $\mathcal{L}_{\mathcal{D}}(\mathbf{f}_*) = \sigma^2$  is the optimal risk achieved by  $\mathbf{f}_* = (\mathbf{B}, \{\mathbf{h}_t\}_{t=1}^T)$ . Let  $\mathcal{D}_{\mathcal{T}}$  be a new distribution with  $y = \mathbf{h}_T^\top \mathbf{B}\mathbf{x} + z$  where  $\mathbf{x}, z$  are independent  $\mathcal{O}(1), \mathcal{O}(\sigma)$  subgaussian respectively and  $\mathbb{E}[\mathbf{x}\mathbf{x}^\top] = \mathbf{I}_p$ . With probability at least  $1 - e^{-cM} - 2e^{-\sqrt{tM}\wedge t}$ , the transfer learning risk on  $\hat{\mathbf{B}}$  with  $M$  samples is bounded as

$$\mathcal{L}_{\mathcal{T}}(\hat{\mathbf{f}}) - \sigma^2 \lesssim r\varepsilon \wedge 1 + C\frac{r+t}{M}.$$

where  $\wedge$  is the minimum symbol. Additionally, if target task vector  $\mathbf{h}_T$  is uniformly drawn from unit Euclidean sphere, in expectation over  $\mathbf{h}_T$  and in probability over target training datasets (with probability at least  $1 - e^{-cM} - 2e^{-\sqrt{tM}\wedge t}$ ), we have the tighter bound

$$\mathbb{E}_{\mathbf{h}_T}[\mathcal{L}_T(\hat{f})] - \sigma^2 \lesssim \varepsilon + C \frac{r+t}{M}.$$

Finally, in both cases, population-empirical transfer gaps  $|\mathcal{L}_T(\hat{f}) - \hat{\mathcal{L}}_{S_T}(\hat{f})|$ ,  $\mathbb{E}_{\mathbf{h}_T}[|\mathcal{L}_T(\hat{f}) - \hat{\mathcal{L}}_{S_T}(\hat{f})|]$  are bounded by  $\mathcal{O}(\frac{r+t}{M})$  with same probability.

**Proof** Let  $\boldsymbol{\theta}_t = \mathbf{B}^\top \mathbf{h}_t$  and  $\hat{\boldsymbol{\theta}}_t = \hat{\mathbf{B}}^\top \hat{\mathbf{h}}_t$ . We first observe that task  $t$  risk is simply

$$\mathcal{L}_t(\hat{\boldsymbol{\theta}}_t) = \mathbb{E}_{\mathcal{D}_t}[(y - \hat{\mathbf{h}}_t^\top \hat{\mathbf{B}} \mathbf{x})^2] = \sigma^2 + \mathbb{E}_{\mathcal{D}_t}[(\boldsymbol{\theta}_t^\top \mathbf{x} - \hat{\boldsymbol{\theta}}_t^\top \mathbf{x})^2] = \sigma^2 + \|\hat{\boldsymbol{\theta}}_t - \boldsymbol{\theta}_t\|^2.$$

Thus, the excess MTL risk is simply

$$\mathcal{L}_{\mathcal{D}}(\hat{\mathbf{f}}) - \mathcal{L}_{\mathcal{D}}(\mathbf{f}_*) = \frac{1}{T} \|\boldsymbol{\Theta} - \hat{\boldsymbol{\Theta}}\|_F^2 \leq \varepsilon,$$

where  $\boldsymbol{\Theta}, \hat{\boldsymbol{\Theta}} \in \mathbb{R}^{T \times p}$  are the concatenated task vectors.

Now, we aim to obtain the transfer learning risk over  $\hat{\mathbf{B}}$ . We first write the target regression task  $(y, \mathbf{x}) \sim \mathcal{D}_T$  with  $y = \mathbf{x}^\top \boldsymbol{\theta}_T + z$  (for some  $\mathbf{h}$ ) as

$$y = \mathbf{x}^\top \hat{\mathbf{B}}^\top \mathbf{h} + z + \mathbf{x}^\top \Pi_{\hat{\mathbf{B}}^\perp}(\boldsymbol{\theta}_T). \quad (47)$$

Here set  $\mathbf{x}' = \hat{\mathbf{B}} \mathbf{x}$  and  $z' = \mathbf{x}'^\top (\mathbf{I} - \hat{\mathbf{B}}^\top \hat{\mathbf{B}}) \boldsymbol{\theta}_T$ , and then  $y = \mathbf{h}^\top \mathbf{x}' + z + z'$ . Note that

$$\mathbb{E}[\mathbf{x}' z'] = \mathbb{E}[\hat{\mathbf{B}} \mathbf{x} \mathbf{x}^\top (\mathbf{I} - \hat{\mathbf{B}}^\top \hat{\mathbf{B}}) \boldsymbol{\theta}_T] = 0,$$

verifying that we can treat the representation mismatch as a dependent but orthogonal subgaussian noise. Combined with Lemma 8, with probability  $1 - e^{-cM} - 2e^{-\sqrt{tM}\wedge t}$  (conditioned on  $\hat{\mathbf{B}}$ ) this leads to a transfer learning risk of

$$\mathcal{L}_T(\hat{f}) - \sigma^2 \leq \|\Pi_{\hat{\mathbf{B}}^\perp}(\boldsymbol{\theta}_T)\|^2 + C \frac{r+t}{M}.$$

Following the proof of Lemma 8, the  $\sigma^2$  term on the right hand side of (42) is related to the inputs ( $\mathbf{x}'$ ) and noise ( $z + z'$ ) levels, which are  $\mathcal{O}(1)$ .

To proceed, observe that  $\mathbb{E}[z'^2] = \|\Pi_{\hat{\mathbf{B}}^\perp}(\boldsymbol{\theta}_T)\|^2 = \|\mathbf{B} \boldsymbol{\theta}_T\|^2 - \|\hat{\mathbf{B}} \boldsymbol{\theta}_T\|^2 = \boldsymbol{\theta}_T^\top (\mathbf{B}^\top \mathbf{B} - \hat{\mathbf{B}}^\top \hat{\mathbf{B}}) \boldsymbol{\theta}_T$ . In the worst case, this risk is equal to

$$\sup_{\|\boldsymbol{\theta}\|=1, \mathbf{B}^\top \mathbf{B} \boldsymbol{\theta} = \boldsymbol{\theta}} \|\Pi_{\hat{\mathbf{B}}^\perp}(\boldsymbol{\theta})\|^2 = \|\mathbf{B}^\top \mathbf{B} - \hat{\mathbf{B}}^\top \hat{\mathbf{B}}\|.$$

Recall that we are given  $\frac{1}{T} \|\boldsymbol{\Theta} - \hat{\boldsymbol{\Theta}}\|^2 \leq \frac{1}{T} \|\boldsymbol{\Theta} - \hat{\boldsymbol{\Theta}}\|_F^2 \leq \varepsilon$ . Additionally,  $\boldsymbol{\Theta}^\top \boldsymbol{\Theta} / T$  is a well-conditioned matrix over the subspace  $\text{Range}(\mathbf{B})$  with minimum nonzero eigenvalue at least  $c/r > 0$  and condition number upper bounded by  $C/c$  (equal to that of  $\mathbf{H}$ ). If  $\varepsilon \leq c/2r$ , this also implies  $\lambda_{\min}(\hat{\boldsymbol{\Theta}}^\top \hat{\boldsymbol{\Theta}} / T) \geq c/2r$  and condition number at most  $3C/c$ . Consequently, applying Davis-Kahan theorem (Yu, Wang, and Samworth 2015) on  $\boldsymbol{\Theta}, \hat{\boldsymbol{\Theta}}$  pair implies that the eigenspaces  $\mathbf{B}, \hat{\mathbf{B}}$  of  $\boldsymbol{\Theta}, \hat{\boldsymbol{\Theta}}$  obey

$$\|\mathbf{B}^\top \mathbf{B} - \hat{\mathbf{B}}^\top \hat{\mathbf{B}}\| \leq \mathcal{O}\left(\frac{\varepsilon}{c/r}\right) = \mathcal{O}(r\varepsilon).$$

If  $r\varepsilon \geq c/2$ , we can simply use the tighter estimate  $\|\mathbf{B}^\top \mathbf{B} - \hat{\mathbf{B}}^\top \hat{\mathbf{B}}\| \leq 1$ , which completes the proof of first part of the lemma.

Secondly, consider the average case scenario where  $\mathbf{h}_T \sim \text{unif\_over\_sphere}$ . In this case, we observe that, the target-averaged transfer risk follows

$$\mathbb{E}_{\boldsymbol{\theta}_T}[\min_{\mathbf{h}} \|\boldsymbol{\theta}_T - \hat{\mathbf{B}}^\top \mathbf{h}\|^2] = \mathbb{E}[\|(\mathbf{I} - \hat{\mathbf{B}}^\top \hat{\mathbf{B}}) \boldsymbol{\theta}_T\|^2] = \frac{1}{2r} \|\mathbf{B}^\top \mathbf{B} - \hat{\mathbf{B}}^\top \hat{\mathbf{B}}\|_F^2. \quad (48)$$

This time, Davis-Kahan theorem yields the tighter estimate (for our purposes)  $\frac{1}{2r} \|\mathbf{B}^\top \mathbf{B} - \hat{\mathbf{B}}^\top \hat{\mathbf{B}}\|_F^2 \lesssim \frac{1}{T} \|\boldsymbol{\Theta} - \hat{\boldsymbol{\Theta}}\|_F^2 \leq \varepsilon$ . To proceed, we find that, the expected transfer learning risk over task distribution obey the tighter guarantee (with probability at least  $1 - e^{-cM} - 2e^{-\sqrt{tM}\wedge t}$ )

$$\mathbb{E}_{\mathbf{h}_T}[\mathcal{L}_T(\hat{f})] - \sigma^2 \lesssim \varepsilon + C \frac{r+t}{M}.$$

The final claim arises as a direct result of our application of Lemma 8 in (47). ■

The following corollary is a Multipath MTL guarantee for least-squares regression obtained by specializing the more general Theorem 8.

**Corollary 5** Suppose  $\mathcal{X} \subset \mathcal{B}^p(R)$ ,  $\ell(\hat{y}, y)$  is quadratic, and Assumptions 3&4 hold. Solving (M<sup>2</sup>TL) with the fixed choice of ground-truth pathways and  $NT \gtrsim \text{DoF}(\mathcal{F}) \log(NT)$ , with probability at least  $1 - \delta$ , we have that

$$\mathcal{R}_{\text{M}^2\text{TL}}(\hat{\mathbf{f}}) \lesssim \sqrt{\frac{L \cdot \text{DoF}(\mathcal{F}) + \log(2/\delta)}{NT}}. \quad (49)$$

**Proof** We need to verify the assumptions of Theorem 8. Observe that  $\mathbf{x}, z$  are subgaussian and ground-truth model  $\|\theta_t^*\| \leq 1$  and all feasible task hypothesis  $\theta$  obeys  $\|\theta\| \leq C^{L+1}$  which we treat as a constant (i.e. fixed depth  $L$ ). Consequently, subexponential norm obeys  $\|(y - \theta^\top \mathbf{x})^2\|_{\psi_1} = \|(z + \mathbf{x}^\top (\theta - \theta_*)^2)\|_{\psi_1} \leq \mathcal{O}(C^{2(L+1)})$  which verifies  $\mathcal{O}(1)$  subexponential condition. Similarly, loss function is Lipschitz with  $\Gamma = \sup_{y, \mathbf{x}} |y - f_t(\mathbf{x})| \leq 2C^{L+1}R$ . Together, these verify Assumption 7. Note that, Theorem 8 has logarithmic dependence on  $\Gamma$  which is subsumed within  $\lesssim$ . Finally, each module is  $C$  Lipschitz (due to spectral norm bounds) and log-covering number of  $d \times p$  matrices with  $C$ -bounded spectral norm obeys  $dp \log(3CR/\varepsilon)$ . These two verify Assumption 8. ■

**Finalizing the Proof of Theorem 4** Following the discussion above, we provide a proof of Theorem 4. The result below is a formal restatement of the theorem with a few caveats. First, we state two closely-related guarantees. First guarantee is when target head  $\mathbf{h}_\mathcal{T}$  is arbitrary (worst-case) and second one is for when it is uniformly distributed over unit sphere (average case). The latter shaves a factor of  $p_L$  in the MTL risk term. Second, the probability term in Theorem 4 is chosen to be approximate for notational simplicity. Namely, we ignored the  $\log(1/\delta)/NT$  term and second order effects. We state the full dependence here which is a bit more convoluted.

**Theorem 9** Suppose Assumptions 3–6 hold and  $\ell(\hat{y}, y) = (y - \hat{y})^2$ . Additionally assume input space is  $\mathcal{B}^p(c\sqrt{p})$  and  $\mathcal{H}_\mathcal{T} = \mathbb{R}^{p_L}$ . Solve MTL problem (M<sup>2</sup>TL) with the knowledge of ground-truth pathways  $(\bar{\alpha}_t)_{t=1}^T$  to obtain a supernet  $\hat{\phi}$  and assume  $NT \gtrsim \text{DoF}(\mathcal{F}) \log(NT)$ . Solve transfer learning problem (TLOP) with  $\hat{\phi}$  to obtain a target hypothesis  $\hat{f}_\phi$ . Then, with probability at least  $1 - 3e^{-cM} - 4\delta$ , excess target risk (3) of TLOP obeys

$$\mathbb{E}_{\alpha_\mathcal{T}}[\mathcal{R}_{\text{TLOP}}(\hat{f}_\phi)] \lesssim \frac{p_L}{M} + p_L \sqrt{\frac{L \cdot \text{DoF}(\mathcal{F}) + \log(2/\delta)}{NT}} + \left[ \sqrt{\frac{\log(2|\mathcal{A}|/\delta)}{M}} \right]_+. \quad (50)$$

Here, the probability is over the source datasets and the (input, noise) pairs of the target dataset i.e.  $(\mathbf{x}_i^T, z_i^T)_{i=1}^M$ , and we used the short hand  $[x]_+ = x + x^2$ . Additionally, if target distribution follows the same generative model with prediction head  $\mathbf{h}_\mathcal{T}$  drawn uniformly at random over the unit sphere, we obtain the tighter bound

$$\mathbb{E}_{\alpha_\mathcal{T}, \mathbf{h}_\mathcal{T}}[\mathcal{R}_{\text{TLOP}}(\hat{f}_\phi)] \lesssim \frac{p_L}{M} + \sqrt{\frac{L \cdot \text{DoF}(\mathcal{F}) + \log(2/\delta)}{NT}} + \left[ \sqrt{\frac{\log(2|\mathcal{A}|/\delta)}{M}} \right]_+. \quad (51)$$

**Remark.** Note that, above, we split probability space into three independent variables. Source datasets  $\mathcal{S}_{\text{all}}$ , (input, noise) pairs of the target dataset i.e.  $(\mathbf{x}_i^T, z_i^T)_{i=1}^M$ , and finally target path  $\alpha_\mathcal{T}$ . The result is with high probability over the former two and expectation over the latter.

**Proof** In this proof, we aim to reduce the Multipath MTL guarantee to a Vanilla MTL scenario so that we can utilize Lemma 9. Assumptions 5 and 6 will be critical towards this goal. Recall that, we have the Multipath MTL guarantee from Corollary 5 so that, with probability  $1 - \delta$ ,

$$\mathcal{R}_{\text{M}^2\text{TL}}(\hat{\mathbf{f}}) \lesssim \sqrt{\frac{L \cdot \text{DoF}(\mathcal{F})}{NT}} + \sqrt{\frac{\log(2/\delta)}{NT}}.$$

Let us call this event  $\mathcal{E}_1$ . Here, we omitted the  $\log |\mathcal{A}|/N$  term because our transfer guarantee will require the knowledge of ground-truth pathways for sources (even if it is not required for the target). The main idea is to show that small  $\mathcal{R}_{\text{M}^2\text{TL}}(\hat{\mathbf{f}})$  implies that target will fall on a pathway with small source-averaged risk. This way, we can apply Lemma 9 to provide a guarantee for the target. To proceed, we gather all unique ground-truth pathways via  $\Gamma = \{\gamma_i\}_{i=1}^S$ . Additionally, let  $C(\gamma)$  be the number of tasks that chooses pathway  $\gamma$ .

Finally let  $\mathcal{L}'_i(\mathbf{f})$  be the excess task-averaged population risk over  $\gamma_i$ , that is  $\mathcal{L}'_i(\mathbf{f}) = \frac{1}{C(\gamma_i)} \sum_{\bar{\alpha}_t = \gamma_i} \{\mathcal{L}_t(f_t) - \sigma^2\}$ . With this definition, we can write MTL excess risk as

$$\frac{1}{T} \sum_{i=1}^S C(\gamma_i) \mathcal{L}'_i(\mathbf{f}) \lesssim \sqrt{\frac{L \cdot \text{DoF}(\mathcal{F})}{NT}} + \sqrt{\frac{\log(2/\delta)}{NT}}.$$

<sup>4</sup>We make this assumption (no norm constraint unlike MTL phase) since during transfer learning, we simply solve least-squares. Thanks to this, we achieve faster rates.

To proceed, we will view each  $\mathcal{L}'_i$  as a vanilla MTL problem over pathway  $\gamma_i$ . Following Assumption 6, we draw the random pathway  $\alpha_{\mathcal{T}}$  of the target task and it is equal to  $\alpha_{\mathcal{T}} = \gamma_i \in \Gamma$ . Note that this event happens with probability  $\mathbb{P}(\alpha_{\mathcal{T}} = \gamma_i) = C(\gamma_i)/T$ . Conditioned on this, let us control the transfer risk.

Note that during TLOP we will search over all pathways  $\alpha \in \mathcal{A}$ . Denote  $\bar{\mathbf{B}}_{\alpha}, \hat{\mathbf{B}}_{\alpha} \in \mathbb{R}^{pL \times p}$  are the ground-truth and empirical weights of the linear model induced by  $\alpha$ . Denote the transfer learning model over  $\hat{\mathbf{B}}_{\alpha}$  via  $\hat{f}_{\alpha}$ . For any choice of  $\alpha$ , applying Lemma 8, we know that empirical-population transfer gap  $\mathcal{L}_{\mathcal{D}_{\mathcal{T}}}(\hat{f}_{\alpha}) - \hat{\mathcal{L}}_{\mathcal{S}_{\mathcal{T}}}(\hat{f}_{\alpha})$  is bounded by  $\mathcal{O}\left(\frac{pL}{M} + \left[\sqrt{\frac{\log(2/\delta)}{M}}\right]_+\right)$  with probability at least  $1 - 2e^{-cM} - 2\delta$ . This is over the input/noise distribution of target samples (arbitrary  $\alpha_{\mathcal{T}} = \gamma_i$  and associated ground-truth  $\theta_{\mathcal{T}}$ ). Union bounding over all potential pathways target task may use, we obtain that,

$$\sup_{\alpha \in \mathcal{A}} |\mathcal{L}_{\mathcal{D}_{\mathcal{T}}}(\hat{f}_{\alpha}) - \hat{\mathcal{L}}_{\mathcal{S}_{\mathcal{T}}}(\hat{f}_{\alpha})| \leq \mathcal{O}\left(\frac{pL}{M} + \left[\sqrt{\frac{\log(2|\mathcal{A}|/\delta)}{M}}\right]_+\right).$$

Consequently, empirical risk minimization over all pathways will choose a target model  $\hat{f}_{\hat{\phi}}$  guaranteeing with probability at least  $1 - 2e^{-cM} - 2\delta$

$$\mathcal{R}_{\text{TLOP}}(\hat{f}_{\hat{\phi}}) \leq \min_{\alpha} \mathcal{R}_{\text{TLOP}}(\hat{f}_{\alpha}) + \mathcal{O}\left(\frac{pL}{M} + \left[\sqrt{\frac{\log(2|\mathcal{A}|/\delta)}{M}}\right]_+\right) \quad (52)$$

$$\leq \mathcal{R}_{\text{TLOP}}(\hat{f}_{\gamma_i}) + \mathcal{O}\left(\frac{pL}{M} + \left[\sqrt{\frac{\log(2|\mathcal{A}|/\delta)}{M}}\right]_+\right). \quad (53)$$

The latter line is reasonable because we know that ground-truth pathway  $\alpha_{\mathcal{T}} = \gamma_i$  is a great candidate for being population minima. Applying Lemma 9 again over the path  $\gamma_i$ , with probability  $1 - e^{-cM} - \delta$ , we obtain

$$\mathcal{R}_{\text{TLOP}}(\hat{f}_{\gamma_i}) \leq pL \mathcal{L}'_i(\mathbf{f}) \wedge 1 + C \frac{pL}{M} + \left[\frac{\log(2/\delta)}{M}\right]_+.$$

Combining with above, with probability  $1 - 3e^{-cM} - 3\delta$ , the ERM solution over all pathways obeys

$$\mathcal{R}_{\text{TLOP}}(\hat{f}_{\hat{\phi}}) \leq pL \mathcal{L}'_i(\mathbf{f}) \wedge 1 + \mathcal{O}\left(\frac{pL}{M} + \left[\sqrt{\frac{\log(2|\mathcal{A}|/\delta)}{M}}\right]_+\right).$$

Note that above holds for worst-case prediction head  $\mathbf{h}_{\mathcal{T}}$ . Additionally, applying Lemma 9 again and assuming  $\mathbf{h}_{\mathcal{T}}$  is generated uniformly over the unit sphere, on the same event, we find

$$\mathbb{E}_{\mathbf{h}_{\mathcal{T}}}[\mathcal{R}_{\text{TLOP}}(\hat{f}_{\hat{\phi}})] \leq \mathcal{L}'_i(\mathbf{f}) + \mathcal{O}\left(\frac{pL}{M} + \left[\sqrt{\frac{\log(|\mathcal{A}|/\delta)}{M}}\right]_+\right). \quad (54)$$

Now, for fixed MTL dataset, taking expectation over  $\alpha_{\mathcal{T}}$ , with same probability over the input/noise distribution

$$\mathbb{E}_{\alpha_{\mathcal{T}}}[\mathcal{R}_{\text{TLOP}}(\hat{f}_{\hat{\phi}})] \leq \sum_{i=1}^S \frac{C(\gamma_i)}{T} pL \mathcal{L}'_i(\mathbf{f}) \wedge 1 + \mathcal{O}\left(\frac{pL}{M} + \left[\sqrt{\frac{\log(2|\mathcal{A}|/\delta)}{M}}\right]_+\right) \quad (55)$$

$$\leq pL \mathcal{R}_{\text{M}^2\text{TL}}(\hat{\mathbf{f}}) \wedge 1 + \mathcal{O}\left(\frac{pL}{M} + \left[\sqrt{\frac{\log(2|\mathcal{A}|/\delta)}{M}}\right]_+\right). \quad (56)$$

Let us call this event  $\mathcal{E}_2$  which is independent of  $\alpha_{\mathcal{T}}$ . Union bounding  $\mathcal{E}_1$  and  $\mathcal{E}_2$  (both independent of  $\alpha_{\mathcal{T}}$ ), we obtain the advertised bound (50). We obtain (51) through same argument following the average-case control (54).  $\blacksquare$

## E Not All Optimal MTL Pathways are Good for Transfer Learning

Ideally we would like to prove Theorem 4 without assuming that MTL phase is solved with the knowledge of ground-truth pathways. While we believe this may be possible under stronger assumptions, here, we discuss why this problem is pretty challenging with a simple example on linear representations.

**Setting:** Suppose we have a single layer linear supernet with  $K = 2$  modules each with size  $2R \times p$ . This corresponds to the Cluster MTL model where we simply wish to group the tasks into two clusters and train vanilla MTL over individual clusters. This simple setting will already highlight the issue.

• **Source tasks:** Consider four groups of tasks  $(\Theta_i)_{i=1}^4$  where  $\Theta_i = (\theta_{ij})_{j=1}^{T/4}$ . We assume that  $T/4$  task vectors from  $\Theta_i$  perfectly span an  $R$  dimensional subspace  $S_i$  (at least  $T \geq 4R$ ). Additionally, set  $(S_i)_{i=1}^4$  to be perfectly orthogonal over  $\mathbb{R}^p$ . Also assume that the tasks are linear and noiseless i.e.  $y_{ij} = \mathbf{x}_{ij}^{\top} \theta_{ij}$ .

**Lemma 10** Suppose representation modules  $\mathbf{B}_1, \mathbf{B}_2 \in \mathbb{R}^{2R \times p}$  are constrained to have orthonormal rows. Define the ground-truth pathways where  $\Theta_1, \Theta_2$  are on pathway 1 and  $\Theta_3, \Theta_4$  are on pathway 2. Now, assume that transfer learning task  $\theta_{\mathcal{T}}$  is drawn uniformly at random from the  $2R$  dimensional subspace of one of these pathways. Assume target is linear & isotropic:  $(\mathbf{x}, y) \sim \mathcal{D}_{\mathcal{T}}$  obeys  $y = \mathbf{x}^\top \theta_{\mathcal{T}} + z$  where  $\mathbb{E}[\mathbf{x}\mathbf{x}^\top] = \mathbf{I}_p$  and  $\mathbf{x}, z$  are orthogonal. Then, regardless of the source sample/task sizes  $N, T$  and target sample size  $M$ , there exists an MTL solution such that, excess transfer risk of final target hypothesis  $\hat{f}_{\mathcal{T}}$  obeys

$$\mathbb{E}_{\theta_{\mathcal{T}}}[\mathcal{R}_{TLOP}(\hat{f}_{\mathcal{T}})] \geq c.$$

for some absolute constant  $c > 0$ . Additionally,  $\mathcal{R}_{TLOP}(\hat{f}_{\mathcal{T}}) \geq 0.5$  almost surely as  $R \rightarrow \infty$ .

**Proof** As the reader might have noticed, the argument is straightforward. Create the following MTL solution: Let  $\mathbf{B}_1$  be an orthonormal basis for  $\Theta_1, \Theta_3$  and let  $\mathbf{B}_2$  be an orthonormal basis for  $\Theta_2, \Theta_4$ . Without losing generality, for  $\mathbf{B}_1$ , let us set it so that first  $R$  rows are assigned to  $\Theta_1$  and last  $R$  assigned to  $\Theta_3$  (same for  $\mathbf{B}_2$ ). Note that, we simply swapped  $\Theta_2$  with  $\Theta_3$  in pathway assignments.

Observe that  $\mathbf{B}_1$  and  $\mathbf{B}_2$  achieves zero MTL risk because they contain all task vectors  $\Theta_i = (\theta_{ij})_{j=1}^{T/4}$  in their range and problems are noiseless. What remains to show is that  $\mathbf{B}_1, \mathbf{B}_2$  assignments are poor choices for the target task drawn from either  $\mathbf{B}_1^*$  induced by  $\Theta_1, \Theta_2$  or  $\mathbf{B}_2^*$  induced by  $\Theta_3, \Theta_4$ . Without losing generality, suppose  $\theta_{\mathcal{T}}$  is drawn from  $\mathbf{B}_1^*$ . Observing  $\theta_{\mathcal{T}}$  lies on the combined range of  $\mathbf{B}_1, \mathbf{B}_2$ , and using properties of linear regression with isotropic features, we bound the target transfer risk via

$$\begin{aligned} \mathcal{L}_{\mathcal{T}}(\hat{f}_{\mathcal{T}}) - \mathbb{E}[z^2] &= \min_{i \in \{1,2\}} \mathcal{L}_{\mathcal{T}}(\mathbf{B}_i^\top \hat{\mathbf{h}}_{\mathcal{T}}) - \mathbb{E}[z^2] \\ &= \min_{i \in \{1,2\}} \|\mathbf{B}_i^\top \hat{\mathbf{h}}_{\mathcal{T}} - \theta_{\mathcal{T}}\|^2 \\ &\geq \min_{i \in \{1,2\}} \min_{\mathbf{h}} \|\mathbf{B}_i^\top \mathbf{h} - \theta_{\mathcal{T}}\|^2 \\ &= \min_{i \in \{1,2\}} \|\mathbf{B}_{3-i} \theta_{\mathcal{T}}\|^2 = \min_{i \in \{1,2\}} \|\mathbf{B}_i \theta_{\mathcal{T}}\|^2 \\ &= \min_{i \in \{1,2\}} \|\text{Proj}_{S_i}(\theta_{\mathcal{T}})\|^2. \end{aligned}$$

The last line highlights the fact that  $S_i$  lies on  $\mathbf{B}_i$  and projection of  $\theta_{\mathcal{T}}$  on  $\mathbf{B}_i$  is exactly equal to its projection on  $S_i$  by pathway assignments. Since  $\theta_{\mathcal{T}}$  is uniformly drawn, the last line is equivalent to  $X(\mathbf{g}, \mathbf{g}') = \frac{\|\mathbf{g}\|^2}{\|\mathbf{g}\|^2 + \|\mathbf{g}'\|^2} \wedge \frac{\|\mathbf{g}'\|^2}{\|\mathbf{g}\|^2 + \|\mathbf{g}'\|^2}$  for  $\mathbf{g}, \mathbf{g}' \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \mathbf{I}_R)$ . Observing  $\|\mathbf{g}\|^2, \|\mathbf{g}'\|^2$  are Chi-squared, it is clear that, for all  $R$  and for some constant  $c_0 > 0$ , we have  $\mathbb{P}(0.5 \leq \frac{\|\mathbf{g}\|^2}{\|\mathbf{g}\|^2 + \|\mathbf{g}'\|^2} \leq 1.5) \geq c_0$ . On these events on  $\mathbf{g}, \mathbf{g}'$ , we have  $X(\mathbf{g}, \mathbf{g}') \geq 1/4$  and  $\mathbb{E}[X(\mathbf{g}, \mathbf{g}')] \geq c = c_0^2/4$ . Finally, as dimension  $R \rightarrow \infty$ , we have  $\|\mathbf{g}\|^2/\|\mathbf{g}'\|^2 \rightarrow 1$  almost surely, which similarly implies  $X(\mathbf{g}, \mathbf{g}') \rightarrow 0.5$ . ■

## F Experimental Details on Section 5

We provide further details on the experiments in Section 5 as well as incorporate additional experiments.

### F.1 Algorithms for Vanilla MTL, Cluster MTL, and Multipath MTL

To facilitate faster and more stable convergence of all three algorithms, we used a conventional approach from nonconvex optimization literature which has also been proposed in the context of linear representation learning (Kong et al. 2020a; Sun et al. 2021; Bouniot et al. 2020; Tripuraneni, Jin, and Jordan 2021). Specifically, linear representation learning with Vanilla MTL has a bilinear form similar to matrix factorization. Thus, first-order method to solve Vanilla MTL benefits from proper initialization of the representation. In our experiments, we use such a two-step procedure:

- **Initialization:** At the start of MTL, build an initialization for the representation.
- **Alternating least-squares (ALS):** Train prediction heads and representation layers through alternating least-squares.

Here, we note that ALS is same as alternating gradient descent (AGD) however we are essentially running infinitely many gradient iterations before alternating. The reason we use this procedure for all three algorithms is to provide a fair comparison without the worry of tuning learning rates for each algorithm individually. Initialization plays a useful role in further stabilizing ALS.

While prior works provide initialization methods for MTL, we will also develop a novel initialization algorithm for Multipath MTL. We believe this may be an interesting future direction for providing provable computational guarantees for Multipath MTL.

**Initialization procedures:** We first revise the procedure for Vanilla MTL. Suppose we are given  $T$  tasks with dataset  $\mathcal{S}_{\text{all}}$  where input features have isotropic covariance. We will use the procedure discussed in (Sun et al. 2021) where the authors claim improvement over (Tripuraneni, Jin, and Jordan 2021; Kong et al. 2020b).

- **Vanilla MTL:** initialization is a method-of-moments procedure as follows:

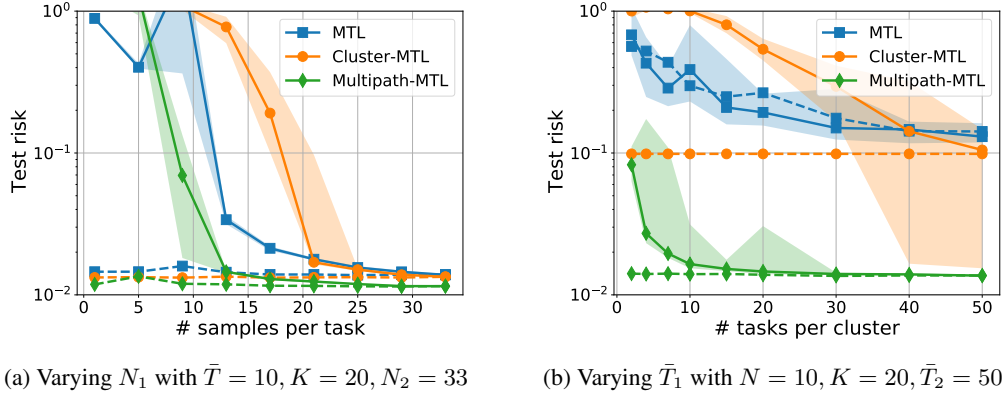


Figure 5: We evaluate (Vanilla) MTL, Cluster-MTL and Multipath-MTL using imbalanced data in a linear regression setting (half tasks have more, half have less data resources). In both experiments, there are  $K = 20$  clusters in total. In Fig. 5a, all clusters have  $\bar{T} = 10$  tasks. Then, we fix the sample size of the tasks in 10 of the clusters ( $N_2 = 33$ ), and change the sample size in the other 10 clusters ( $N_1$ ) from 1 to 33. Then we solve the three MTL (Vanilla, Cluster, Multipath) problems. Solid/Dashed curves show the test risk of the tasks who have fewer/more samples. In Fig. 5b, instead we fix the sample size of each task ( $N = 10$ ). In 10 of the clusters, we set the number of tasks  $\bar{T}_2 = 50$ . While in the other clusters, the number of tasks is varied from 2 to 50. Again, we run experiments under all the three settings and plot the test risk of fewer/more tasks in solid/dashed curves. The curves in both figures show the median risks and the shaded regions highlight the first and third quantile risks. Each marker is average of 20 independent runs.

1. Form the  $\hat{\theta}_t$  estimates via  $\hat{\theta}_t = \frac{1}{N} \sum_{i=1}^N y_{ti} \mathbf{x}_{ti}$ .
2. Form the moment matrix  $M = \sum_{t=1}^T \hat{\theta}_t \hat{\theta}_t^\top$ .
3. Set  $\hat{B}_0 \in \mathbb{R}^{R \times p}$  to be the top  $R$  eigenvectors of  $M$ .

At this point, we can start running our favorite choice of first order method starting from the initialization  $\hat{B}_0$ . In our implementation, we run ALS where we estimate  $\{\hat{h}_t\}_{t=1}^T$  (by fitting LS given  $\hat{B}$ ), then re-estimate  $\hat{B}$  and keep going.

• **Cluster MTL:** In our experiments, we assumed clusters (i.e. pathways) are known. This is in order to decouple the challenge of task-clustering from the comparisons in Figure 2. We note that task clustering has been studied by (Fifty et al. 2021; Kumar and Daume III 2012; Kang, Grauman, and Sha 2011) (Leveraging relations between tasks are explored even more broadly (Zhuang et al. 2020), however these works don't come with comparable statistical guarantees. In our setup, Cluster MTL simply runs  $K$  Vanilla MTL algorithms in parallel over individual clusters using ground-truth pathways.

• **Multipath MTL:** We propose an initialization algorithm which is inspired from the Vanilla MTL algorithm as follows. Again, we assume knowledge of clustering/pathways.

1. Estimate shared first layer  $\hat{B}_1 \in \mathbb{R}^{R \times p}$  via **Vanilla MTL** initialization using all data.
2. Estimate cluster-specific representations  $(\tilde{B}_2^k)_{k=1}^K \in \mathbb{R}^{r \times p}$  via **Vanilla MTL** initialization over each cluster data.
3. Estimate the second layer  $(\hat{B}_2^k)_{k=1}^K \in \mathbb{R}^{r \times R}$  by projecting  $\tilde{B}_2^k$  onto the  $R$ -dimensional first layer as follows

$$\hat{B}_2^k = \tilde{B}_2^k \hat{B}_1^\top.$$

We then run ALS where we go in the order: Prediction heads, second layers, first layer (repeat).

**Remark on unknown clusters:** We note that a simple approach to identifying clusters when they are unknown is by solving Vanilla MTL and then clustering the resulting weight vectors  $\{\hat{\theta}_t\}_{t=1}^T$  of the Vanilla MTL solution (e.g. via  $K$ -means). The reason is that, the ground-truth weights  $\{\theta_t^*\}_{t=1}^T$  are simply points that lie on  $r$ -dimensional latent cluster-subspaces that we would like to recover. Naturally, the (random) points on the same subspace will have higher correlation. This viewpoint (restricted to linear setting) also connects well with the broader subspace clustering literature where each learning task is a point on a high-dimensional subspace (Vidal 2011; Parsons, Haque, and Liu 2004; Elhamifar and Vidal 2013). The challenge in our setting is we only get to see the points through the associated datasets. Figure 3 shows our results assuming unknown source pathways.

In the next section, we discuss a few more experiments comparing these three approaches.

## F.2 Additional Numerical Experiments

In Figure 5, we conduct more experiments to see how tasks with less data resources perform in MTL when trained together with other tasks which have more resources. Here, by resources we either mean a task having more samples  $N$  or a task having other

(related) tasks along its pathway/cluster. Thus, our experiments involve imbalanced training data. We consider two experimental settings to show how Multipath MTL benefits accuracy compared to the other two MTL models: Vanilla MTL and Cluster MTL. **Experimental settings:** Consider the same Vanilla MTL, Cluster MTL and Multipath MTL problems in linear regression regime as discussed in Section 5 and follow the same algorithm in Section F.1. In the experiments, same as Section 5, we set  $p = 32$ ,  $R = 8$ , and  $r = 2$ . We consider MTL problem with  $K = 20$  clusters. Here, data is noisy. In Fig. 5a, there are 10 tasks in each cluster. In half of the clusters, each task has fixed sample size,  $N_2 = 33$  (more resource); while in the remaining 10 clusters, the sample size ( $N_1$ ) varies from 1 to 33 (less resource). Solid curves display the test risk of the tasks with  $N_1$  samples and dashed curves present the test risk of tasks with  $N_2$  samples. Rather than changing number of samples, in experiments shown in Figure 5b, we create another scenario where number of tasks per cluster is varied (as a measure of data resource). Here, instead all tasks contain  $N = 10$  samples. For 10 of the total clusters, there are fixed  $\bar{T}_2 = 50$  tasks in each cluster. However, the other 10 contain only  $\bar{T}_1$  tasks in each cluster, and we compare the performance with different  $\bar{T}_1$  selections. We change  $\bar{T}_1$  from 2 to 50 and results are displayed in Fig 5b. Similar, solid curves present the results of the clusters who contain fewer tasks (less resource), to the contrary, dashed curves present the test risk of clusters with fixed  $\bar{T}_2 = 50$  tasks (more resource).

In both figures, Multipath-MTL performs better than the other two models, which again shows that the sample complexity of hierarchical model is smaller than the vanilla and clustering models. When there are fewer samples or fewer tasks, all the three methods fail at learning a good representation. The three dashed curves in Fig. 5a behave in line with expectations: They follow from the fact that tasks with more samples can learn decent representations by themselves. The solid curve of Cluster MTL decreases slower, and it is because other than the other two methods where clusters are correlated and representations are shared, in Cluster MTL setting (as depicted in Fig. 4b), clusters are separately trained. Therefore, there is no benefit across the clusters. In Fig. 5b, firstly, the evidence that orange and blue dashed curves are above the green one again shows the sample efficiency of Multipath MTL. Here, when there are only 2 tasks for the 10 resource-poor clusters, the Cluster MTL has the worst performance because there is no representation sharing across clusters. Test risk of Vanilla MTL does not change too much even the task number increases. It is because MTL representation of vanilla model is larger and tasks don't have enough samples to train their prediction heads. For instance, blue solid curve hits blue dashed curve at very beginning, which shows that the model is already trained well and adding more tasks cannot help too much (both more resource tasks and less resource tasks are doing similar).