#### CANCER

# DORGE: Discovery of Oncogenes and tumoR suppressor genes using Genetic and Epigenetic features

Jie Lyu<sup>1</sup>\*, Jingyi Jessica Li<sup>2</sup>\*<sup>†</sup>, Jianzhong Su<sup>3</sup>, Fanglue Peng<sup>3</sup>, Yiling Elaine Chen<sup>2</sup>, Xinzhou Ge<sup>2</sup>, Wei Li<sup>1†</sup>

Data-driven discovery of cancer driver genes, including tumor suppressor genes (TSGs) and oncogenes (OGs), is imperative for cancer prevention, diagnosis, and treatment. Although epigenetic alterations are important for tumor initiation and progression, most known driver genes were identified based on genetic alterations alone. Here, we developed an algorithm, DORGE (Discovery of Oncogenes and tumor suppressoR genes using Genetic and Epigenetic features), to identify TSGs and OGs by integrating comprehensive genetic and epigenetic data. DORGE identified histone modifications as strong predictors for TSGs, and it found missense mutations, super enhancers, and methylation differences as strong predictors for OGs. We extensively validated DORGE-predicted cancer driver genes using independent functional genomics data. We also found that DORGE-predicted dual-functional genes (both TSGs and OGs) are enriched at hubs in protein-protein interaction and drug-gene networks. Overall, our study has deepened the understanding of epigenetic mechanisms in tumorigenesis and revealed previously undetected cancer driver genes.

Copyright © 2020 The Authors, some rights reserved; exclusive licensee American Association for the Advancement of Science. No claim to original U.S. Government Works. Distributed under a Creative Commons Attribution NonCommercial License 4.0 (CC BY-NC).

#### INTRODUCTION

Cancer results from an accumulation of key genetic alterations that disrupt the balance between cell division and apoptosis (1). Genes with "driver" mutations that affect cancer progression are known as cancer driver genes (2), which can be classified as tumor suppressor genes (TSGs) and oncogenes (OGs) based on their roles in cancer progression (3). OGs are usually activated by gain-of-function mutations that stimulate cell growth and division, whereas TSGs are inactivated by loss-of-function (LoF) mutations (frameshift insertions/deletions and nonsense mutations) that block TSG functions in inhibiting cell proliferation, promoting DNA repair, and activating cell cycle checkpoints.

CRISPR-Cas9 screens with libraries of single-guide RNAs are powerful tools for identifying genes essential for cancer cell fitness, such as cancer cell growth and viability. For example, recent CRISPR screens by the Wellcome Sanger Institute detected 628 priority targets in 324 human cell lines from 30 cancer types (4). However, the genes identified by CRISPR screens in cell lines, which differ vastly from primary cells, may not be physiologically relevant to human biology and disease. Many well-known cancer driver genes in the Cancer Gene Census (CGC) database (5) were missing in CRISPR-screening results. They might have phenotypic effects in animal models that are not included in the current CRISPR screens.

Hence, it is necessary to predict cancer driver genes based on patient genomics data. Cancer genome sequencing efforts, such as The Cancer Genome Atlas (TCGA) (6), have generated an unprecedentedly large data resource and enabled the development of bioinformatics algorithms to discover cancer driver genes. Tokheim *et al.* (7) reviewed eight major algorithms, and Bailey *et al.* (8) integrated 26 computational tools in a pan-cancer mutation study. These

algorithms mainly look for cancer driver genes with greater than expected background mutational rates, and they output a ranked list of candidate genes based on a small collection of genetic features such as somatic mutations and copy number alterations (CNAs). Tumor Suppressor and OG Explorer (TUSON) (9) and the 20/20+ machine-learning method (7) are the two major algorithms that can distinguish between protein-coding TSGs and OGs based on distinct patterns of mutational signatures.

However, a recent meta-analysis indicated that, over the next 10 years, even if all available tumor genomes were analyzed, many cancer driver genes would remain undetected because of the lack of distinction between driver mutations and background mutational load (10). In addition, emerging evidence suggests that genetic alterations alone are insufficient to explain all cancer driver genes, including some well-known ones. For example, sustained expression of estrogen receptor- $\alpha$  (ESR1) drives two-thirds of breast cancers, but ESR1 mutations that alter transcription levels occur in only 7% of ESR1-positive tumors (11). Furthermore, many pediatric tumors have extremely low mutation rates; some even appear to have no substantial recurrent somatic mutations (12). Thus, it is likely that other mechanisms, such as epigenetic alterations, are responsible for the dysregulation of many cancer driver genes.

For example, trimethylation on histone H3 lysine 4 (H3K4me3) and DNA methylation are the most extensively studied epigenetic modifications that influence gene expression and cell fate. H3K4me3 is a widely recognized marker of active promoters and regulates the preinitiation complex formation and gene activation (13). More than 80% of promoters containing H3K4me3 are transcribed (14), and H3K4me3 is also involved in pre-mRNA splicing, recombination, DNA repair, and enhancer function. DNA methylation occurs in 70 to 80% of 5′—C—phosphate—G—3′ (CpG) sites in a normal genome (15). H3K4me3 and CpG methylation alteration are associated with disease initiation, including many types of cancer (16). In particular, promoter hypermethylation that silences TSGs is a key epigenetic event in tumorigenesis (17), whereas gene-body methylation is positively correlated with gene expression (18). Recently, the "broad epigenetic domain" has emerged as a new concept in the control of

<sup>&</sup>lt;sup>1</sup>Division of Computational Biomedicine, Department of Biological Chemistry, School of Medicine, University of California, Irvine, Irvine, CA 92697, USA. <sup>2</sup>Department of Statistics, University of California, Los Angeles, Los Angeles, CA 90095, USA. <sup>3</sup>Department of Molecular and Cellular Biology, Baylor College of Medicine, Houston, TX 77030, USA.

<sup>\*</sup>These authors contributed equally to this work.

<sup>†</sup>Corresponding author. Email: wei.li@uci.edu (W.L.); jli@stat.ucla.edu (J.J.L.)

cancer development. In an integrative analysis of 1134 genome-wide ChIP-seq datasets (19) from the Encyclopedia of DNA elements (ENCODE) project (20), we found that broad H3K4me3 is a unique epigenetic signature of TSGs. In contrast to the common sharp (e.g., <1-kb width) H3K4me3 peaks associated with increased transcriptional initiation, broad H3K4me3 peaks are associated with increased transcriptional elongation. In addition, we also found many wide gene-body regions that are lowly methylated in normal tissues (the regions called "gene-body methylation canyons") as hypermethylated in cancer (21). Gene-body methylation canyons are unexpectedly enriched in OGs, and their hypermethylation directly induces OG activation (21).

Nevertheless, to the best of our knowledge, none of the existing bioinformatics algorithms sufficiently leveraged epigenetic features to predict cancer driver genes, despite the fact that epigenetic alterations are known to be associated with cancer driver genes. Therefore, these algorithms were not fully empowered, and there is a pressing need for a computational algorithm that integrates epigenetic data with genetic alterations to improve the prediction of cancer driver genes.

To address this need, we developed DORGE (Discovery of Oncogenes and tumor suppressoR genes using Genetic and Epigenetic features). DORGE includes two prediction algorithms: DORGE-TSG for predicting TSGs and DORGE-OG for predicting OGs; both algorithms are elastic net-based logistic regression (LR) classifiers trained on CGC genes and neutral genes (NGs). By evaluating DORGE-TSG and DORGE-OG, we found an unusually large contribution of histone modification to TSG prediction, as well as crucial roles of the features such as missense mutations, genomics, super enhancer percentages, and hypermethylation in predicting OGs. Cancer driver genes predicted by DORGE include known cancer driver genes and novel ones that have not been reported in the literature. We evaluated these novel cancer driver genes using multiple genomics and functional genomics datasets. In addition, we found that the novel dual-functional genes, which DORGE predicted as both TSGs and OGs, are highly enriched at hubs in proteinprotein interaction (PPI) and drug/compound-gene networks.

#### **RESULTS**

### DORGE predicts TSGs and OGs based on known cancer driver genes and NGs

We developed a computational tool DORGE, by integrating extensive genomic and epigenomic datasets, for predicting cancer driver genes, i.e., TSGs and OGs. Briefly, we used CGC genes and 75 curated candidate features to train two binary classification algorithms: DORGE-TSG and DORGE-OG, which we subsequently applied to every gene to predict its probability of being a TSG and OG, respectively. Last, we used the predicted probabilities to rank genes genome-wide and identified the top-ranked genes as candidate TSGs and OGs.

Prediction of cancer driver genes is a classification problem. It requires a high-quality training dataset that contains reliable TSGs, OGs, and the genes unlikely to be TSGs or OGs. Our two positive-training gene sets include 242 TSGs and 240 OGs (with dual-functional genes removed) from the CGC database v.87, which we refer to as CGC-TSGs and CGC-OGs hereafter. The negative-training gene set includes 4058 NGs reported to have no cancer relevance (9). To allow for the prediction of dual-functional genes that are

both a TSG and an OG, we trained two classifiers for predicting TSGs and OGs, respectively.

To develop DORGE, we constructed 75 features that are likely predictive of cancer driver genes based on the literature. These features have either known roles in TSG/OG disruption (e.g., DNA methylation and somatic mutations) or potential links to TSG/OG functions (e.g., CRISPR-screening data; data file S1). We categorized these features into four major types: (i) 33 mutational features from two well-known cancer driver gene prediction algorithms— TUSON (9) and 20/20+ (7)—and Genome Aggregation Database (gnomAD); 28 of these 33 features were compiled by TCGA (6) and Catalogue of Somatic Mutations in Cancer (COSMIC) (5) from the mutation data of patient samples; (ii) 12 genomic features including 3 from 20/20+ (7) and 9 features (e.g., gene lengths and genome evolution-related features) that have not been previously used to predict cancer driver genes (22); (iii) 27 epigenetic features, including histone modifications from the ENCODE project (20), promoter and gene-body methylation features from the COSMIC database, and super enhancer percentages from the dbSUPER database (23); and (iv) 3 phenotypic features, including CRISPR-screening data from the DepMap project (24), Variant Effect Scoring Tool (VEST) scores from 20/20+(7), and gene expression Z scores from TCGA.

To train classifiers for TSG and OG prediction, we compared eight classification algorithms: LR, LR with the lasso penalty, LR with the ridge penalty, LR with the elastic net penalty, random forests, support vector machines (SVM) with the linear kernel, SVM with the Gaussian kernel, and XGBoost (https://github.com/ dmlc/xgboost). For each algorithm, we considered three class ratios (where a class ratio was defined as the number of NGs to the number of CGC-TSGs or CGC-OGs): the original ratio, 5:1, and 1:1; for the latter two ratios, we randomly divided NGs into partitions so that the number of NGs in each partition approximately met the ratio given the number of CGC-TSGs or CGC-OGs. Considering the imbalance between NGs and CGC-TSGs/CGC-OGs in sizes, we used the fivefold cross-validated (CV) area under the precisionrecall curve (AUPRC), instead of the receiver operating characteristic curve, as the accuracy measure to compare these eight classification algorithms under the three class ratios. Our comparison result showed that downsampling the NGs to have more balanced class ratios as 5:1 and 1:1 did not improve the accuracy achieved by the original class ratio. Hence, we decided to keep the original class ratio and found that LR with the lasso, LR with the ridge, LR with the elastic net, and random forests performed the best with similar AUPRC values (data file S2). We chose LR with the elastic net as the classification algorithm for its good interpretability and its capacity for selecting correlated, informative features. Then, we trained LR with the elastic net separately for TSG and OG prediction and subsequently used the two trained algorithms to assign every gene a TSG score and an OG score, both ranging from 0 to 1, with a larger value indicating a higher chance of the corresponding gene being a TSG or an OG. To decide appropriate thresholds on the TSG scores and OG scores for final predictions, we weighted the severity of mispredicting NGs as TSGs/OGs (i.e., making false-positive predictions) versus the other way around and set a target false-positive rate (FPR) of 1%. Last, we used the Neyman-Pearson classification algorithm (25) to set thresholds on the TSG scores and OG scores by respecting our target FPR and obtained two classifiers: DORGE-TSG and DORGE-OG for predicting TSGs and OGs, respectively.

Next, we identified the important features for TSG and OG prediction. Because many features are correlated (data file S1), the feature coefficients estimated by LR with the elastic net are not biologically interpretable measures of feature importance. The reason is that if one adds to the training data a feature that is highly correlated with an existing feature, the estimated coefficient of the existing feature would become less significant. This phenomenon contradicts our biological interpretation of feature importance: If a feature is important, its importance should not be diluted by the addition of another feature. Yet, we are still interested in the importance of features in our final multifeature linear classifier, so marginal feature importance based on each feature alone does not suffice. To address this issue, we proposed a simple two-step procedure: (i) we clustered features into feature groups that were approximately uncorrelated with one another, and (ii) we evaluated the importance of each feature group by the reduction in the fivefold CV AUPRC when that feature group was left out, i.e., the contribution of that feature group to the fivefold CV AUPRC given all the other feature groups. Our simple but innovative approach is advantageous in three aspects. First, by grouping correlated features, we can interpret a small number of feature groups, each of which has a distinct biological interpretation, instead of a large number of features. Second, making feature groups approximately uncorrelated has a desirable consequence: if a new feature were added, it would either be added to an existing feature group or create a new feature group by itself (if it is approximately uncorrelated with any existing features); then, its addition would barely affect the importance of the feature groups it is not in, as uncorrelated features would not affect each other's importance in a multifeature linear classifier. Third, the same criterion, fivefold CV AUPRC, was used to select a classification algorithm and define the importance of a feature group, making the analysis self-consistent. Using this approach, we first divided all 75 features into 20 feature groups by hierarchical clustering with complete linkage so that features within each group have pairwise absolute Pearson correlations of at least 0.1 (data file S2). Then, we ranked the 20 feature groups by their contributions to fivefold CV AUPRC and selected the top-ranked groups as those whose contributions exceeded 0.005. This gave us three and five feature groups for predicting TSGs and OGs, respectively.

Analyzing these top predictive feature groups, we found that multiple histone modification features stood out as the most predictive group (whose contribution to 5-fold CV AUPRC was almost 10-fold of that of the second most predictive group containing phenotype features) for TSGs and that missense mutations constituted the top feature group for predicting OGs (Fig. 1, A and B). Besides, epigenetic features including super enhancer and cancer-normal methylation differences in promoter and gene-body regions were among the top feature groups for predicting OGs (Fig. 1B). We also found histone modifications and missense mutations among the top predictive features for both TSGs and OGs (Fig. 1, A and B), suggesting that TSGs and OGs share certain features, whose predictive power for TSGs and OGs may be different though. For each feature within a top-ranked TSG (or OG) feature group, we compared its values in the CGC-TSGs (or CGC-OGs) and the NGs by the two-sided Wilcoxon rank-sum test, and the resulting  $-\log_{10}P$ value was shown in Fig. 1 (A and B).

We further examined several features in terms of their individual, marginal power of distinguishing CGC-TSGs and CGC-OGs from NGs. Multiple features are marginally strong predictors of TSGs, as they have significantly higher values in CGC-TSGs than in NGs. They include epigenetic features such as H3K4me3 peak length and height (Fig. 1C and fig. S1A) and H3K79me2 peak length and height (fig. S1, B and C), missense mutational features such as nonsilent/ silent ratio (fig. S1D), and phenotype features such as VEST score (Fig. 1D). Many features also have significantly higher values in CGC-OGs than in NGs. They include missense mutational features such as missense damaging/benign ratio (Fig. 1E), missense entropy (Fig. 1F), probability of being LoF intolerant (pLI) score (Fig. 1G), and LoF o/e (observed/expected) constraint (fig. S1E); genomics features such as evolutionary conservation phastCons score and noncoding Genomic Evolutionary Rate Profiling (ncGERP) score (fig. S1, F and G); and epigenetic features such as super enhancer percentage in cell lines (Fig. 1H). In particular, our finding agrees with previous studies in that missense damaging/benign ratio (reflecting the functional impact of missense mutations) and missense entropy (representing the enrichment of mutations in few residues) (9) have significantly higher values in CGC-OGs than in CGC-TSGs and NGs (Fig. 1, E and F). VEST and PolyPhen-2 scores, both of which reflect functional effects of mutations, have significantly higher values in CGC-TSGs and CGC-OGs than in NGs, and they do not exhibit statistically significant differences between CGC-TSGs and CGC-OGs (Fig. 1D and fig. S1H). We found super enhancer, a commonly regarded OG-specific feature (26), also characteristic of TSGs, as it has significantly higher values in CGC-TSGs than in NGs (Fig. 1H).

We note that, besides H3K4me3 peak length, a readily known TSG predictor, peak lengths of four more histone marks (H3K79me2, H3K36me3, H4K20me1, and H3K9ac) are also significantly larger in CGC-TSGs than in CGC-OGs and NGs (fig. S1, B, I, J, and K), consistent with the fact that the activation of TSGs is associated with transcriptional elongation (19, 27–29). To further verify the enrichment of broad H3K4me3 peaks in CGC-TSGs, we performed the Fisher's exact test on a two-by-two contingency table, whose two rows correspond to CGC-TSGs and all the other genes in the training data (CGC-OGs and NGs) and whose two columns correspond to the genes with broad H3K4me3 peaks (whose mean lengths across ENCODE samples are >4 kb) and the rest of genes. We similarly performed two more Fisher's exact tests to check the enrichment of broad H3K4me3 peaks in CGC-OGs and NGs but found much lower enrichment in these two gene groups than in CGC-TSGs, confirming that H3K4me3 is a distinctive feature of TSGs (fig. S1L). Together, we identified histone modifications as the top predictors for TSGs. We found missense mutations, super enhancer percentages, and methylation differences between cancer and normal samples as major predictors for OGs. It is worth noting that histone modifications and missense mutations are also important features for predicting OGs and TSGs, respectively, although to a lesser extent. In summary, DORGE can successfully leverage public data to discover the genetic and epigenetic alterations that play significant roles in cancer driver gene dysregulation. Figure S2 provides an overview of the DORGE method and the evaluations in the following sections.

#### **Evaluation of the prediction accuracy of DORGE**

As we described earlier, DORGE-TSG and DORGE-OG output TSG scores and OG scores for predicting TSGs and OGs, respectively. Every gene received a TSG score and an OG score, both ranging from 0 to 1, and a higher TSG score (or OG score) indicates a

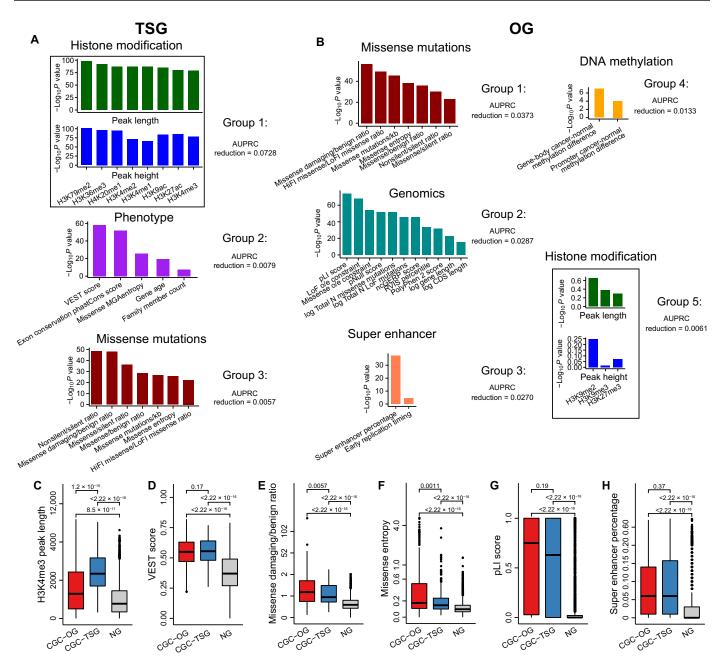


Fig. 1. Features that discriminate TSGs from OGs. (A) Feature groups selected for TSGs. (B) Feature groups selected for OGs. Feature groups are sorted according to the AUPRC reduction in elastic net fivefold cross-validation. Feature groups are named according to the representative features. Box plots showing the distribution of (C) H3K4me3 mean peak length, (D) VEST score, (E) missense damaging/benign ratio, (F) missense entropy, (G) pLI score, and (H) super enhancer percentage for the CGC-OG, CGC-TSG, and NG sets. Genes as both TSGs and OGs are excluded. P values for the differences between the TSGs/OGs and NGs were calculated by the one-sided "greater than" Wilcoxon rank-sum test.

higher probability of a gene being a TSG (or an OG; Materials and Methods). DORGE thresholded the TSG scores and OG scores by the Neyman-Pearson classification algorithm (25) with a target FPR of 1%, leading to 925 predicted TSGs, whose TSG scores exceeded 0.6233374, and 683 predicted OGs, whose OG scores exceeded 0.6761319. In total, DORGE predicted 1172 cancer driver genes, including 436 dual-functional genes (Fig. 2A; the predicted genes are listed in data file S2). We note that these predicted TSGs and OGs

are conservative predictions guided by the small FPR threshold 1%, as reflected by the fact that their numbers are smaller than the numbers of previously predicted cancer driver genes—1217 TSGs and 803 OGs in databases TSGene (30) and ONGene (31) (by 18 June 2020). If DORGE users would like to be less conservative and predict more TSGs and OGs, they can opt for a higher FPR threshold such as 5%. Next, we filtered out CGC genes from the DORGE-predicted cancer driver genes and defined the remaining 725 predicted

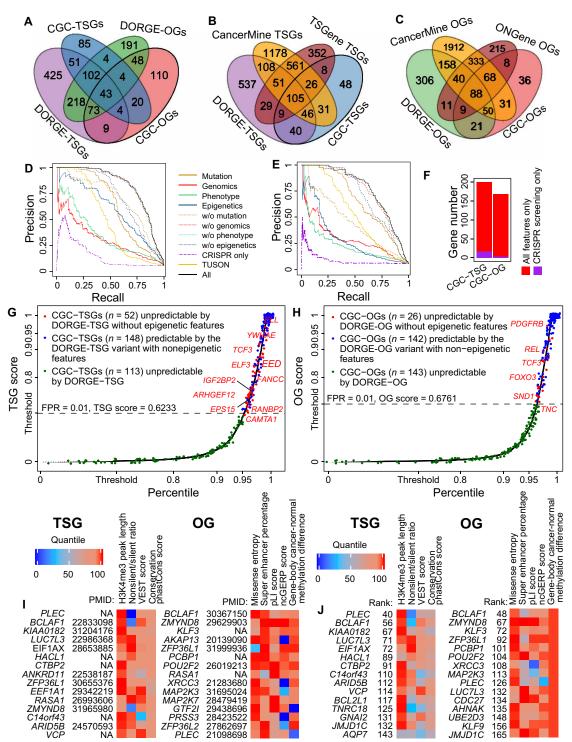


Fig. 2. Evaluation of the DORGE method and characterization of the DORGE-predicted novel TSGs and OGs. Venn diagrams showing the overlap (A) between DORGE-predicted novel TSGs/OGs and CGC-TSGs/OGs; (B) between DORGE-predicted novel TSGs, CGC-TSGs, CancerMine-TSGs, and TSGene database-TSGs; and (C) between DORGE-predicted novel OGs, CGC-OGs, CancerMine-OGs, and ONGene database-OGs. Precision-recall curves (PRCs) for (D) TSG and (E) OG prediction. Different lines represent different PRCs from DORGE or DORGE variants. (F) Stacked bar plots showing the number of rediscovered CGC-TSGs and CGC-OGs using all features compared to CRISPR-screening data only. Cumulative distribution function (CDF) plots of DORGE-predicted TSG scores (G) and OG scores (H) of 19,636 human genes. The *x* axis and the *y* axis are swapped for illustration purposes, and the *y* axis is stretched to emphasize large TSG and OG scores. CGC genes are plotted as jitter points to avoid overplotting. The dashed lines indicate DORGE-TSG and DORGE-OG thresholds at a target FPR of 1%, and the CGC genes whose TSG scores and OG scores exceed the thresholds (above the dashed lines) are predicted as TSGs and OGs. (I) Top 15 DORGE-predicted non-CGC novel TSGs (left) and OGs (right), respectively, along with representative feature heatmaps and PubMed IDs. To make features comparable, feature values are transformed into quantiles. (J) Top 15 DORGE-predicted non-CGC novel TSGs (left) and OGs (right) that have no documented role in cancer based on the TSGene, ONGene, and CancerMine databases, along with representative feature heatmaps.

TSGs and 515 predicted OGs as DORGE-predicted novel genes (data file S1), among which 537 novel TSGs were not included in the CancerMine (32) or TSGene database (Fig. 2B), and 306 novel OGs were not found in the CancerMine or ONGene database (Fig. 2C).

We evaluated DORGE-TSG and DORGE-OG by their overall prediction accuracy and found that they achieved high fivefold CV AUPRC of 0.821 and 0.766, respectively, when trained with all the 75 features (Fig. 2, D and E). Considering that previous algorithms primarily relied on genetic features to predict cancer driver genes, we evaluated the accuracy gain of DORGE from including epigenetic and phenotypic features. To this end, we constructed variants of DORGE-TSG and DORGE-OG based on each of the following feature subsets: "Mutation," "Genomics," "Phenotype," "Epigenetics," and their complements (i.e., the subsets resulting from subtracting each of the four feature subsets from the 75 features), as well as TUSON and CRISPR screening—only features (Fig. 2, D and E, and data file S1). For each of these DORGE-TSG and DORGE-OG variants, we calculated its fivefold CV AUPRC.

On the basis of feature subsets Mutation, Genomics, Phenotype, and Epigenetics, the corresponding DORGE-TSG variants achieved fivefold CV AUPRC of 0.638, 0.314, 0.358, and 0.600, respectively. In parallel, on the basis of the complements of Mutation, Genomics, Phenotype, and Epigenetics (i.e., when features in each subset were excluded), the corresponding DORGE-TSG variants achieved fivefold CV AUPRC of 0.692, 0.819, 0.820, or 0.715. These results consistently show the large contributions of Mutation and Epigenetics features to TSG prediction (Fig. 2D). Furthermore, using the TUSON method's features and the CRISPR screening-only feature, the corresponding DORGE-TSG variants only achieved fivefold CV AUPRC of 0.500 and 0.156, much lower than 0.821 achieved by DORGE-TSG with all the 75 features. Similarly, we compared DORGE-OG with its variants trained on feature subsets. Specifically, DORGE-OG variants that only used Mutation, Genomics, Phenotype, or Epigenetics features achieved fivefold CV AUPRC of 0.660, 0.299, 0.241, or 0.295, respectively; when each of these feature subsets was excluded, the AUPRC correspondingly became 0.453, 0.752, 0.763, or 0.705. These results suggest that Mutation features have a large contribution to OG prediction (Fig. 2E). Similar to DORGE-TSG, the DORGE-OG variants trained with TUSON features or the CRISPR screening-only feature had much lower prediction accuracy (fivefold CV AUPRC of 0.534 or 0.089) than that of DORGE-OG trained with all the 75 features (fivefold CV AUPRC of 0.766). The fact that DORGE-TSG and DORGE-OG outperformed all their variants confirms that DORGE effectively leveraged the 75 features and did not suffer from overfitting in its TSG and OG prediction.

The above results also reveal that the CRISPR screening-only feature did not have a high predictive power on its own, as shown by its low fivefold CV AUPRC (0.156 and 0.089) in TSG and OG prediction. Moreover, under the target FPR of 1%, the DORGE-TSG and DORGE-OG variants with the CRISPR screening-only feature identified only 16 (5.1%) CGC-TSGs and 3 (1.0%) CGC-OGs, whereas DORGE-TSG and DORGE-OG with all the 75 features recovered an additional 184 (58.8%) CGC-TSGs and 165 (53.1%) CGC-OGs (Fig. 2F). These results challenge a common belief that CRISPR screening using cell lines is powerful for discovering cancer driver genes. A possible reason for our results is that cell lines do not mimic in vivo cancer cells well. These additional cancer driver

genes with all the 75 features might have phenotypic effects in animal models that have not been included in the current CRISPR screens.

We next evaluated the distinct predictive power provided by epigenetic features to cancer driver gene prediction. Inspecting the distributions of TSG scores and OG scores, we found that many top-ranked CGC genes were not predictable by DORGE without epigenetic features (Fig. 2, G and H). In detail, 52 (16.61%) CGC-TSGs and 26 (8.36%) CGC-OGs would have been missed by DORGE-TSG and DORGE-OG, respectively, at the target FPR of 1% if epigenetic features were not included. These results suggest that (i) epigenetic features empowered the discovery of cancer driver genes and (ii) epigenetic features empowered DORGE-TSG more than DORGE-OG because the number of rescued CGC-TSGs (52) is twice the number of rescued CGC-OGs (26).

We then searched biomedical literature for the top 15 novel TSGs and OGs ranked by DORGE. Out of these top novel genes, 10 TSGs and 12 OGs have reported tumor-suppressive and oncogenic functions, respectively (Fig. 2I). We also inspected these top novel genes for selected representative features and confirmed that they have high values in the top predictive TSG features (H3K4me3 peak length, nonsilent/silent ratio, VEST score, and conservation phastCons score) and OG features (missense entropy, super enhancer percentage, pLI score, ncGERP score, and gene-body cancernormal methylation difference) selected from the top feature groups (Fig. 2I). We further confirmed this result in the subset of top novel genes that are not in the CancerMine, TSGene, and ONGene databases (Fig. 2J). In particular, nearly all of the top novel TSGs have broad H3K4me3 peaks, and most of the top novel OGs are hypermethylated in gene body (with positive cancer-normal methylation differences).

#### Benchmarking DORGE against existing algorithms

We further compared DORGE with 10 existing algorithms for cancer driver gene prediction using four accuracy measures—sensitivity (Sn), specificity (Sp), precision, and overall accuracy—all based on CGC genes (Table 1). We did not include the five-test model (RF5) because although it outputs TSG and OG probabilities, it does not have explicit cutoffs for defining TSGs and OGs (33). We found that DORGE performed the best in all these measures except Sp, for which DORGE was 0.997 and the best algorithm 20/20+ was 1.000. The superiority of DORGE was most obvious in Sn, where its top performance (0.611) was followed with a large gap by OncodriveFM (0.338) (34), MuSIC (0.331) (35), and MutPanning (0.318) (36) (Table 1). To further confirm that DORGE outperformed these 10 algorithms, we performed a similar comparison based on 1056 OncoKB cancer genes (37), which had been widely used to benchmark cancer gene prediction. Consistent with the CGC gene evaluation results, DORGE achieved the best performance in Sn (almost 50% higher than that of the second best algorithm OncodriveFM) and overall accuracy, the third best performance in Sp (0.997 versus 0.999 of the best method TUSON), and the second best performance in precision (0.973 versus 0.993 of the best method 20/20+, whose *Sn* was only 32% of that of DORGE) (data file S2). Together, our benchmark results show that DORGE made a significant advance in improving cancer driver gene prediction from existing algorithms.

On the basis of CGC-TSGs and CGC-OGs, we further benchmarked DORGE against 20/20+, TUSON, and Genes Under Selection in Tumors (GUST) for separate prediction of TSGs and OGs

Method	#	Sn	Sp	Precision	Accuracy	Algorithms
DORGE	1172	0.611	0.997	0.966	0.948	Logistic regression with the elastic net model
OncodriveFM (34)	2600	0.338	0.915	0.367	0.841	Functional impact model
MuSIC (35)	1975	0.331	0.870	0.272	0.801	Mutational background model
MutPanning (36)	460	0.318	0.994	0.880	0.907	Nucleotide context model
TUSON (9)	243	0.222	0.999	0.961	0.900	<i>P</i> value combination
OncodriveFML (58)	680	0.212	0.983	0.646	0.885	Functional impact model
20/20+ (7)	193	0.208	1.000	0.991	0.899	Random Forest model
GUST (78)	276	0.206	0.994	0.838	0.894	Random Forest model
MutSigCV (57)	158	0.137	0.998	0.905	0.888	Mutational background model
OncodriveCLUST (59)	586	0.118	0.963	0.319	0.855	Mutational hotspot model
ActiveDriver (61)	417	0.098	0.996	0.771	0.881	Logistic regression model

(data file S2). We did not include the other seven algorithms because they could not predict TSGs and OGs separately. Consistent with our previous results, DORGE exhibited much higher Sn than the other three algorithms (DORGE had Sn of 0.639 and 0.54 for predicting TSGs and OGs, while the best Sn of the other three algorithms was only 0.252 and 0.116), and it also achieved the best precision and overall accuracy; all the four algorithms had close to perfect Sp. Although the high Sn of DORGE seemed to be due to the fact that 20/20+, TUSON, and GUST by default predicted fewer TSGs and OGs than DORGE did, it was not the case. After we adjusted the thresholds of 20/20+ and TUSON so that they predicted the same numbers of TSGs and OGs as DORGE did (the GUST software does not allow such threshold adjustment), the Sn of 20/20+ and TUSON, though increased, remained almost onefold lower than that of DORGE. Collectively, our results suggest that DORGE outperformed 20/20+, TUSON, and GUST in both TSG and OG prediction.

We also compared DORGE with TUSON and 20/20+ in terms of their predicted ranking of CGC-TSGs and CGC-OGs. For example, if an algorithm predicted gene A more likely than gene B to be a TSG, we say that gene A received a higher TSG rank (smaller in rank number) than gene B did. Accordingly, we calculated a TSG rank and an OG rank for every CGC gene by each algorithm. Among the CGC genes, we define the core CGC-TSGs and core CGC-OGs as those that were annotated solely as TSGs and OGs, not both (dualfunctional), in CGC v.77. Compared to the genes that were added later to CGC v.87, these core CGCs have been more extensively studied. Then, we examined the ranking consistency between DORGE and the other two algorithms for CGC genes and the core CGC genes. For CGC-TSGs, we found that their TSG ranks by DORGE had strong positive correlations with their TSG ranks by

TUSON and 20/20+ (fig. S3, A and B), and overall, they were ranked higher by DORGE than by the other two algorithms (fig. S3E). We observed similar results for CGC-OGs (fig. S3, C, D, and G). The conclusions also held for core CGC genes (fig. S3, F and H). These results confirm that DORGE predictions are more biologically relevant than those of TUSON and 20/20+. For example, ELL (elongation factor for RNA polymerase II), a CGC-TSG, was ranked 190th by DORGE-TSG, 8144th by TUSON, and 3958th by 20/20+; PDGFB (platelet-derived growth factor subunit B), a CGC-OG, was ranked 207th by DORGE, 2753rd by TUSON, and 4982nd by 20/20+. Also, DORGE ranked CGC dual-functional genes better than TUSON and 20/20+ did, as exemplified by the dual-functional gene IDH1 {isocitrate dehydrogenase [nicotinamide adenine dinucleotide phosphate (NADP+)] 1}, which was ranked first for TSG and 28th for OG by DORGE, 18,734th for TSG and 2092nd for OG by TUSON, and 14,936th for TSG and 13th for OG by 20/20+.

### Functional evaluation of novel cancer driver genes and those unpredictable without epigenetics features

Even though DORGE predicted many more cancer driver genes than TUSON, 20/20+, and GUST did—DORGE, TUSON, 20/20+, and GUST predicted 1172, 243, 193, and 276 cancer driver genes, respectively—DORGE achieved the highest overall prediction accuracy based on CGC genes. After confirming this, we further characterized the novel cancer driver genes, defined as those predicted by DORGE but not included in the CGC database.

We performed the Kyoto Encyclopedia of Genes and Genomes (KEGG) pathway analysis on the novel TSGs and OGs, and we found, as expected, that the novel TSGs are enriched with TSG-related pathways such as "apoptosis" and "focal adhesion" and that the

novel OGs are enriched with OG-related pathways such as "cell cycle" and "transforming growth factor– $\beta$  (TGF- $\beta$ ) signaling pathway" (Fig. 3A). However, without epigenetic features, the novel TSGs and OGs predicted by the DORGE-TSG and DORGE-OG variants are no longer enriched with certain TSG-related and OG-related pathways such as TGF- $\beta$  signaling pathway (fig. S4A). These results again suggest that epigenetic features made unique contributions to discovering novel cancer driver genes. In addition, the degrees of enrichment ( $-\log_{10} P$  values) of those shared enriched KEGG pathways, which were enriched in novel TSGs or OGs regardless of the inclusion of epigenetic features, are positively correlated, implying that the addition of epigenetic features did not prohibit the discovery of meaningful cancer driver genes (fig. S4, B and C).

Given that histone modification features (e.g., H3K4me3 peak length) empowered DORGE-TSG prediction, we sought experimental evidence for the novel TSGs that have broad histone modification (e.g., H3K4me3) peaks. A previous cell proliferation experiment observed increased cell growth after knocking down multiple potential TSGs whose H3K4me3 peaks have mean lengths (across ENCODE cell lines) greater than 2 kb (19), including two DORGE-predicted novel TSGs, CSRNP1 and NR3C1B. Another previous study found that Mll4 loss down-regulates potential TSG expression and weakens broad H3K4me3 peaks in mice (38). Examining the human orthologs of the six mouse potential TSGs that were down-regulated by Mll4 loss in that study, we found that four orthologs were ranked top by DORGE-TSG and have H3K4me3 peaks longer than 2 kb. These four human genes are DNMT3A (18th), BCL6 (96th), FOXO3 (222nd), and CBFA2T3 (1012th).

### Characterization of DORGE-predicted novel TSGs and OGs by independent functional genomics data

We first used a published Transposase Accessible Chromatin with high-throughput sequencing (ATAC-seq) dataset of TCGA pancancer samples (39) to characterize the DORGE-predicted novel cancer driver genes. ATAC-seq reveals gene accessibility and provides valuable information about the complex gene regulatory relationships. On the basis of this ATAC-seq dataset, we found that DORGE-predicted novel TSGs and OGs, consistent with CGC-TSGs and CGC-OGs, are significantly more accessible than NGs (all with  $P=2.22\times10^{-16}$  by the one-sided Wilcoxon rank-sum test; Fig. 3B). This result established a connection between cancer driver genes and chromatin accessibility; both TSGs and OGs are ubiquitously accessible in cancer samples.

We then explored a possible relationship between cancer driver genes and epigenetic regulators (ERs), which are known to play fundamental roles in genome-wide gene regulation by reading or modifying chromatin states. A previous study suggested that most ERs are intolerant to LoF mutations (40), and our fig. S1E also shows that LoF mutations (reflected by the LoF o/e constraint feature) are significantly more abundant in TSGs and OGs than NGs, prompting us to explore whether ER genes have a significant overlap with cancer driver genes. By analyzing a curated list of 761 ERs, we found significant enrichment of CGC-TSGs and CGC-OGs  $(P = 3.14 \times 10^{-20})$  and  $P = 9.36 \times 10^{-8}$ , respectively, by the Fisher's exact test; in total, 94 CGC cancer driver genes are among the ERs, with  $P = 2.79 \times 10^{-13}$  by the Fisher's exact test; Fig. 3C). This result also shows the greater enrichment of CGC-TSGs than that of CGC-OGs in ER genes, consistent with a previous study showing that the application of cancer gene classifiers to ER genes revealed more

TSGs than OGs (41). Similar to CGC genes, DORGE-predicted novel TSGs ( $P = 1.15 \times 10^{-6}$ ) are also more enriched than novel OGs ( $P = 2.65 \times 10^{-3}$ ) in ER genes (Fig. 3C).

We next evaluated DORGE-predicted novel TSGs using Sleeping Beauty (SB) screening data. The SB transposon is a type of synthetic DNA element that can disrupt the expression of genes near its insertion sites, a process called insertional mutagenesis. Hence, the SB transposon is a screening tool for TSGs, whose expression disruption leads to carcinogenesis. To verify the novel TSGs, we downloaded the list of inactivating pattern genes from the SB Cancer Driver Database (SBCDDB) (42). As expected, we found that both CGC-TSGs ( $P = 5.41 \times 10^{-19}$ ) and DORGE-predicted novel TSGs ( $P = 5.11 \times 10^{-24}$ ) are enriched in the list. In contrast, NGs have no enrichment. This result is consistent with our expectation that TSGs are inactivated in SB screens (Fig. 3D).

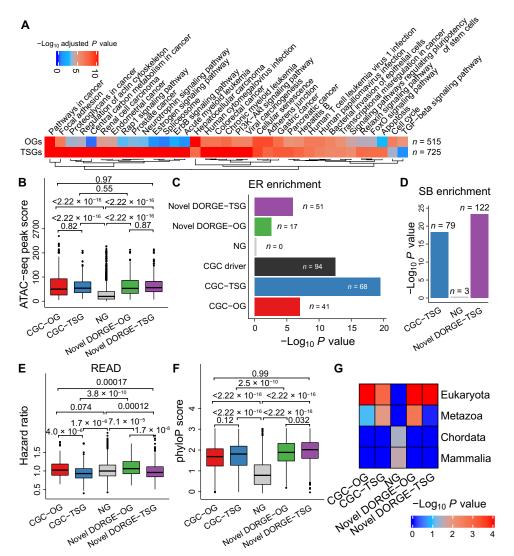
We further evaluated DORGE-predicted novel cancer driver genes using an shRNA screening dataset from the Achilles project (43), as shRNA screens for gene essentiality for cell proliferation in cell lines. On the basis of the dataset, the knockdown of DORGE-predicted novel OGs and CGC-OGs shows a greater decrease in cell proliferation rates compared to NGs (fig. S4D). In contrast, the knockdown of DORGE-predicted novel TSGs and CGC-TSGs shows nearly no decrease in cell proliferation rates compared to NGs (fig. S4D). This result is consistent with the prior knowledge that the proliferation of cell lines is dependent upon OGs (24).

Last, we evaluated DORGE-predicted novel cancer driver genes using patient survival data. In the precomputed survival data downloaded from the OncoRank website (44), every gene has a hazard ratio (HR; whose value >, =, or <1 indicates that the gene's expression reduces, does not affect, or increases patients' survival time, respectively). We found that CGC-TSGs and DORGE-predicted novel TSGs have significantly lower HRs than OGs (CGC-OGs and DORGE-predicted novel OGs) and NGs in three representative cancer types: rectum adenocarcinoma, colon adenocarcinoma, and uterine corpus endometrial carcinoma (Fig. 3E and fig. S4, E and F). These results are consistent with the fact that TSG expression prohibits cancer occurrence and prolongs survival, while OG expression has the opposite effects. The complete HRs and P values of DORGE-predicted novel TSGs and OGs in 21 cancer types are available in data file S1.

### TSGs and OGs are conserved at both exons and noncoding regions

Previous studies have suggested that evolutionarily conserved genes are enriched with cancer driver candidates and drug targets (45). Consistent with these studies, we observed statistically significant differences in exonic sequence conservation (phastCons and phyloP scores) between CGC-TSGs/OGs and NGs, and the same conclusion holds for DORGE-predicted TSGs and OGs (Fig. 3F and fig. S4G). Compared to OGs, TSGs have slightly higher exonic sequence conservation (Fig. 3F and fig. S4G).

We next explored the conservation of noncoding regions in cancer driver genes. Noncoding regions are characterized by positive ncGERP values and negative noncoding Residual Variation Intolerance Score (ncRVIS) values. The reason is that ncGERP is a measure of nucleotide constraints and reflects conservation across the mammalian lineage (fig. S1G) (46), while ncRVIS measures humanspecific constraints (46). On the basis of these two measures, we found that TSGs (CGC-TSGs and DORGE-predicted novel



**Fig. 3.** Characterization and evaluation of DORGE-predicted novel TSGs/OGs by independent functional genomic and genomic datasets. (A) KEGG pathway enrichment analysis performed by Enrichr (75) for DORGE-predicted novel TSGs and OGs. Because of space limitations, terms with adjusted *P* values <10<sup>-4</sup> are shown. Besides, terms with adjusted *P* values 10<sup>8</sup>-fold lower for TSGs than OGs or 10<sup>4</sup>-fold lower for OGs than TSGs are also shown. (**B**) ATAC-seq peak score measuring open chromatin for CGC-TSGs/OGs, DORGE-predicted novel TSGs/OGs, and NGs. Enrichment heatmaps of various gene types in (**C**) ER gene list and (**D**) inactivating pattern gene list for SB insertional mutagenesis, a screening tool for cancer driver genes. (**E**) Boxplot showing the Cox hazard ratio (HR) score for various gene types. Data are from rectum adenocarcinoma (READ). (**F**) Boxplot showing the phyloP score for various gene types. The phyloP score measures phylogenetic conservation and represents – log*P* values under a null hypothesis of neutral evolution. PhyloP basewise conservation scores were derived from a Multiz alignment of 46 vertebrate species. (**G**) TSGs and OGs are enriched in genes having earlier evolutionary origin (Eukaryota). *P* values for the differences between indicated gene categories were calculated by the one-sided Wilcoxon rank-sum test. In boxplots and heatmap, the Fisher's exact test is used to calculate *P* values, and gene numbers in different gene categories are normalized to 200 to make *P* values comparable. In this figure, dual-functional CGC genes were excluded from the CGC-TSGs/OGs.

TSGs) are slightly more conserved than OGs (CGC-OGs and DORGE-predicted novel OGs) at noncoding regions (figs. S1G and S4H).

In summary, we found that cancer driver genes are more conserved than NGs at both exonic and noncoding regions. Between TSGs and OGs, we found that TSGs are more conserved at exons, while OGs are more conserved at non-coding regions.

#### TSGs and OGs are overrepresented in ancient genes

Motivated by our conservation results, we investigated the phyletic ages (i.e., evolutionary origins) of cancer driver genes. Although

cancer driver genes are believed to have originated from Metazoa (multicellular animals) (47), the possibility of their origination from Eukaryota, an earlier evolutionary origin, has not been explicitly investigated. On the basis of the phyletic-age gene lists (from early to late: Eukaryota, Metazoa, Chordata, and Mammalia) from the Online GEne Essentiality (OGEE) database (48), we found significant enrichment of cancer driver genes in the Eukaryota gene list (Fig. 3G; P values by the Fisher's exact test:  $P = 1.05 \times 10^{-3}$  for CGC-TSGs,  $P = 3.25 \times 10^{-13}$  for DORGE-predicted novel TSGs,  $P = 1.41 \times 10^{-5}$  for CGC-OGs, and  $P = 2.77 \times 10^{-5}$  for DORGE-predicted novel OGs), in contrast to NGs. Our results indicate that

cancer driver genes may have originated earlier in the evolutionary history than previously thought. In addition, we found that cancer driver genes were not enriched in young phyletic ages (Chordata and Mammalia; Fig. 3G), consistent with a recent paper (49).

### Dual-functional cancer driver genes act as backbones in PPI networks

Previous studies have shown high interactivity of cancer driver genes in the BioGRID PPI network (9), and accordingly, PPI data have been used to identify cancer driver genes (50, 51). We, therefore, explored the extent to which DORGE-predicted TSGs and OGs are connected to other genes/proteins. When analyzing the whole BioGRID PPI network (Fig. 4A), we found that TSGs and OGs, including CGC genes and DORGE-predicted novel genes, exhibit significantly higher degrees, betweenness, and closeness centrality than NGs do (fig. S5, A to C). This result suggested that the removal or knockdown of cancer driver genes, as expected, will exert a critical impact on the whole PPI network, in particular, dual-functional driver genes, as both TSGs and OGs, display even higher interactivity than sole TSGs and OGs do (fig. S5, A to C). Densely connected genes tend to form modules, and cancer driver gene modules can trigger the hallmarks of cancer and confer the proliferation advantages displayed on cancer cells (52). Here, we used the Molecular Complex Detection (MCODE) algorithm to identify six densely connected network modules/backbones (Fig. 4B) from the PPI subnetwork of the 1172 DORGE-predicted cancer driver genes. The 64 genes that comprise the six identified modules are all dual-functional genes (8 CGC dual-functional genes and 56 DORGE-predicted novel dual-functional genes). This overrepresentation of dual-functional driver genes in network modules is unusual, as it is highly unlikely to obtain a 64-gene subnetwork composed of all dual-functional genes ( $P = 6.66 \times 10^{-27}$  by the bino-

It was previously shown that somatic alterations often occur at PPI network hub genes in cancer (53), and these hub genes are typically essential genes. We therefore investigated the enrichment of cancer driver genes in the hub genes, the 978 genes (top 5%) with the highest degrees in the BioGRID PPI network. We found that all TSGs, OGs, and dual-functional genes (including CGC genes and DORGE-predicted novel genes) are enriched in the hub genes (Fig. 4C). The CGC and novel dual-functional genes are the most enriched (Fig. 4C). We also analyzed the enrichment of 10 functional gene sets. Among these gene sets, we found that the genes with high missense o/e constraints (highest top 5%), the essential genes from the OGEE database, and the ER genes are most enriched in the hub genes (Fig. 4D). Previous literature has not reported any connection between ERs and PPI hub genes, and our finding strengthens the critical roles of ERs. We also found that the genes with broad H3K4me3 peaks are significantly enriched, to a similar degree to the housekeeping genes (HKGs), in the hub genes (Fig. 4D).

#### ER genes act as backbones in gene-drug networks

Cancer driver gene prediction is the basis for the development of anticancer drugs and personalized cancer treatments. We therefore explored possible gene-drug relationships of DORGE-predicted cancer driver genes using the PharmacoDB, a gene-drug network constructed from comprehensive high-throughput cancer pharmacogenomic datasets. In the subnetwork containing CGC genes and DORGE-predicted novel genes, we found that these cancer driver

genes are densely connected to anticancer drugs (fig. S5D). Similar to our observation from the PPI network, we found that TSGs and OGs, including CGC genes and DORGE-predicted novel genes, exhibit significantly denser connections to drugs than NGs do (fig. S5E).

We then identified the top 10 drugs with the largest numbers of connected genes in the PharmacoDB gene-drug network. Among these 10 drugs, the top one is doxorubicin, a well-known chemotherapeutic agent, and the other nine drugs are also known anticancer drugs (fig. S5F). We next identified 979 genes (top 5%) with the highest degrees in the gene-drug network as hub genes and found that DORGE-predicted novel driver genes are enriched in these hub genes (Fig. 4E). We also analyzed the enrichment of 10 functional gene sets in these hub genes. Unlike their enrichment in our previously defined PPI network hub genes (Fig. 4D), the essential genes and the HKGs are not enriched in these gene-drug network hub genes (Fig. 4F), an expected result as their expression is required for normal cells and they are unlikely to be viable drug targets for cancer treatment. In contrast, we still observed the enrichment of three functional gene sets—the genes with high missense o/e constraints (highest top 5%), the ER genes, and the genes with broad H3K4me3 peaks-in the gene-drug network hub genes (Fig. 4F). Together with our PPI analysis, we conclude that the genes in these three functional gene sets may be potential actionable drug targets. To the best of our knowledge, there has been no report on the enrichment of the ER genes in gene-drug network hub genes. Our results from PPI and gene-drug network analysis emphasize the importance of studying the ER genes as potential drug targets.

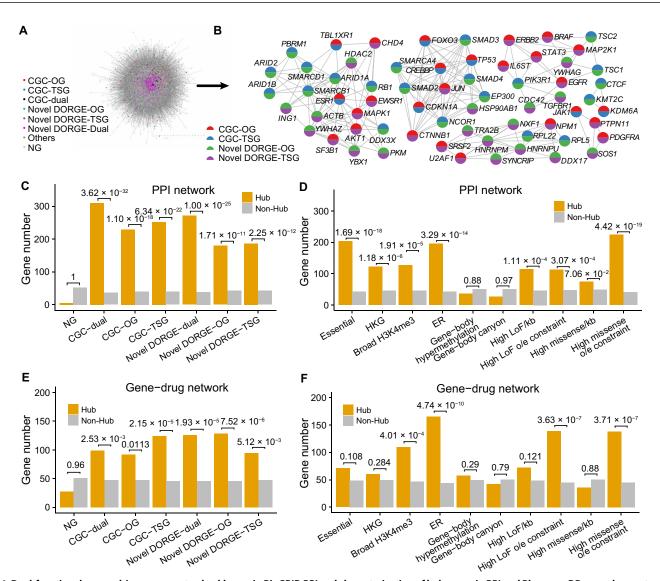
### Identification of candidate anticancer drugs from public transcriptomic data

A bottleneck in novel anticancer drug discovery is an efficient selection of potential molecular targets for a drug/compound or its derivatives. Ideal anticancer drugs are those that up-regulate TSGs and/or down-regulate OGs. We used the CRowd Extracted Expression of Differential Signatures (CREEDS) data (54) to explore the relationship between CGC and DORGE-predicted genes and anticancer drugs (data file S1). We identified 68 proven or potential anticancer drugs/compounds that were associated with 68 target genes meeting the filtering criteria (limma Q value < 0.05 and fold change > 2) from the CREEDS data (fig. S6). Fifty-four (79.41%) of the 68 genes are DORGE-predicted novel TSG or OG genes.

Recent pharmacological efforts suggested that drugs/compounds actionable toward more than one gene or molecular pathway are preferable for repurposing (55), and it is common for existing drugs to be later repurposed as anticancer drugs. For example, dexamethasone was previously classified as a corticosteroid but later repurposed for cancer treatment. Among the 68 drugs/compounds we identified, 30 are anticancer and chemotherapy drugs (fig. S6, bottom), 23 have only been tested in laboratories and are not yet in clinical trials, and 15 have not been tested in cell lines (fig. S6, bottom). Of the 38 drugs/compounds not yet confirmed in anticancer clinical trials, many have been proven to treat other diseases. Overall, our results indicate that they are potential drugs for cancer treatment.

#### **DISCUSSION**

In this paper, we developed a machine-learning tool DORGE for identifying cancer driver genes by integrating genetic and epigenetic features. Our development is the first effort that goes beyond the



**Fig. 4. Dual-functional cancer driver genes act as backbones in BioGRID PPI and characterization of hub genes in PPI and PharmacoDB gene-drug networks.**(A) Complete BioGRID PPI network. (B) The Molecular Complex Detection (MCODE) algorithm was applied to DORGE-predicted novel TSGs/OGs to identify densely connected network modules (or backbones). All genes in the identified network are CGC dual-functional genes or novel dual-functional genes. Gene categories are represented as pie charts, with the colors coded based on gene categories. (C) Enrichment of CGC-TSGs/OGs and DORGE-predicted novel TSGs/OGs in hub genes in the BioGRID network. (D) Enrichment of various gene sets or epigenetic and mutational patterns in hub genes in the BioGRID network. (E) Enrichment of CGC-TSGs/OGs and DORGE-predicted novel TSGs/OGs in hub genes in the PharmacoDB gene-drug network. (F) Enrichment of various gene sets or epigenetic and mutational features in hub genes in the PharmacoDB gene-drug network. Hub genes are defined as the genes with the top 5% highest degree in the BioGRID or PharmacoDB network. To generate comparable *P* values, the gene numbers in different gene categories were normalized to 200. Broad H3K4me3: Genes with H3K4me3 length > 4000. *P* values for the differences between indicated gene categories were calculated by the right-sided Wilcoxon rank-sum test.

use of tumor genetic alterations for cancer driver prediction, and it was motivated by our previous studies that found specific epigenetic patterns associated with TSGs or OGs (19, 21). Although experimental validation is needed for further studies, our computational evaluation verifies that the novel cancer driver genes predicted by DORGE resemble known cancer drivers in multiple aspects and show promise as potential therapeutic targets. In particular, the topranked novel cancer driver genes, especially those regulated by epigenetic mechanisms, warrant further detailed investigation.

Cancer driver genes that are infrequently mutated in cancer are often indistinguishable from passenger genes with random mutations in genome sequencing data. Such random mutations may result from technical reasons including tumor DNA contamination, sequencing depth, and mutation calling failure (56). Therefore, infrequently mutated cancer driver genes are hardly detectable by the methods based on the mutational background model [MutSigCV (57)] or the functional impact model [OncodriveFML (58), OncodriveFM (34), and OncodriveCLUST (59)]. However, these genes may be identified through the integration of epigenetic, phenotypic, and genomic data.

In previous studies, various nonmutational datasets have been used in cancer driver gene identification; however, unlike DORGE,

existing work only used few or several nonmutational features extracted from these datasets (7, 50, 51, 57, 60, 61). For example, MutSigCV used DNA replication timing and cell line expression data (57); ActiveDriver used phosphorylation site information (61); and 20/20+ used multispecies conservation, mutation pathogenicity scores, and replication timing (7). PPI networks and pathway knowledge have also been used to identify cancer driver genes (50, 51); however, these studies were biased toward well-studied genes/pathways and thus may overlook quite many genuine cancer driver genes. In contrast to all these studies, DORGE leverages epigenetic information without any bias toward gene selection to predict cancer driver genes, and this innovation results in DORGE overpowering these existing studies in discovering novel cancer driver genes.

We further note that the capacity of DORGE in predicting TSGs and OGs separately allows DORGE to identify novel dual-functional cancer driver genes. This is advantageous given that more and more dual-functional cancer driver genes have been identified in the literature. In this study, we found a unique property of dual-functional cancer driver genes: They have more connecting partners in PPI and drug-gene networks than other driver genes have. This property, to our knowledge, was not previously reported. Several novel dual-functional genes predicted by DORGE drew our attention. For example, PTEN (phosphatase and tensin homolog), a protein phosphatase, is commonly regarded as a TSG; however, DORGE predicted it as an OG as well. We found that oncogenic roles were reported for PTEN in a few studies. One explanation for the dualfunctional roles of PTEN is that its oncogenic effect depends on the positions of mutations (62). We confirmed this by analyzing the mutation patterns of PTEN and found one pattern as the classic OG mutation pattern with most substitutions in p.R130 (63). In DORGE, further updates can quantify the dual-functional roles (i.e., the relative chance of being TSGs or OGs) of dual-functional genes.

While we have already found dozens of nonmutational features that contribute significantly to the predictive power of DORGE, many CGC genes remain undetected by DORGE (Fig. 2, G and H). A possible reason is the missingness of other factors or mechanisms that regulate cancer driver genes. Fortunately, the continuing increase in functional genetic and epigenetic data will provide a lasting opportunity to improve and fine-tune cancer driver gene prediction methods. In future studies, we can perform lineage-specific rather than pan-cancer prediction and extend DORGE to predicting long noncoding genes, as many features used in DORGE are not restricted to protein-coding genes. In addition, further work is needed for a better understanding of phenomena such as ancient phyletic ages of cancer driver genes and enrichment of cancer driver genes at PPI and gene-drug network hubs.

In summary, this study highlights the integration of epigenetic data to achieve a more comprehensive prediction of cancer driver genes. DORGE will serve as an essential resource for cancer biology, particularly in the development of targeted therapeutics and personalized medicine for cancer treatment.

#### **MATERIALS AND METHODS**

#### **Experimental design**

In this paper, we propose DORGE, a machine-learning framework incorporating a large number of features to discover TSGs/OGs (fig. S2). First, we used CGC v.87 genes and NGs as the training genes to predict TSGs and OGs separately from 75 candidate features

by LR with the elastic net penalty, and the resulting two classifiers are DORGE-TSG and DORGE-OG. Next, we used fivefold cross-validation to evaluate DORGE. We also analyzed the benefit of introducing epigenetic features based on KEGG enrichment and evaluated DORGE based on several genomic and functional genomic datasets. Last, we showed the enrichment of dual-functional genes predicted by DORGE in hub genes in PPI and gene-compound networks.

#### **Gene annotation**

All gene annotations, genomic, and functional genomic datasets were downloaded from hg19 genome version or processed to hg19 if they were from other genome versions. Genome version conversion was done using the LiftOver program (https://genome.ucsc.edu/cgi-bin/hgLiftOver). HUGO Gene Nomenclature Committee annotation (www.genenames.org/) was used for gene annotation. The gene annotation can be found in data file S1. Promoters were defined as the regions from the upstream 1000 base pairs (bp) to downstream 500 bp of transcription start sites (TSSs), while genebody regions were defined as the regions from downstream 500 bp of TSSs to transcription termination sites (TTSs).

### Datasets used in this study Somatic mutation datasets

The somatic mutation dataset used in this study was derived from the TCGA (6) website (https://portal.gdc.cancer.gov/) and COSMIC, v86 (5). These two datasets were combined to help increase the mutational information of infrequently mutated genes. Duplicate tumor samples present in more than one dataset were excluded. The final dataset used for the calculation of mutation-related features contained 5,700,484 mutations from more than 30 tumor types. Hypermutated tumor samples with >2000 mutations were excluded from this dataset. The population genetic dataset (pLoF Metrics by Gene TSV file) for evaluating features, such as LoF intolerance, was downloaded from gnomAD v2 (https://gnomad.broadinstitute.org/downloads/) (64). Additional details regarding features calculation can be found in data file S1.

#### **Epigenetic datasets**

We downloaded all peak BED files (hg19) for H3K4me3 and other representative histone modifications from the ENCODE project (www.encodeproject.org/). The full file names and download links are listed in data file S1. The gene-body canyon annotation file (65), including DNA methylation information, was obtained from a previous study (21), which were based on TCGA whole-genome bisulfite sequencing (WGBS) data. The data for calculating promoter and gene-body cancer-normal methylation difference were also downloaded from the level 3 methylation data from the COSMIC website (v.90). Repli-seq Binary Alignment Map (BAM) datasets were downloaded from the ENCODE project website, and the featureCounts program (http://subread.sourceforge.net/) was used to assign BAM reads to gene bodies. Read counts were normalized on the basis of the sequencing depth of the BAM files, and the normalized read numbers were used to calculate the replication timing S50 score (66). This score, which determines the median replication timing, was calculated by a tool available from a previous study (66). The super enhancer annotation was downloaded from the dbSUPER database (23).

#### Other datasets

The level 3 TCGA data, which include the processed somatic CNA and gene expression data, were downloaded from the COSMIC

website (v.90) and used without processing. The processed cell proliferation (dependency) scores from 436 CRISPR-treated cell line samples were obtained from the DepMap website (Avana-17Q2-Public\_v2) (24). For each gene, gene expression was aggregated across samples to obtain the median Z score. The phastCons scores were downloaded from the University of California, Santa Cruz (UCSC) (http://hgdownload.cse.ucsc.edu/goldenPath/hg19/phastCons46way/). The dataset including the gene damage index (GDI), Primate dN/dS score, RVIS percentile, ncRVIS, ncGERP, family member count, and gene age features were downloaded from https://github.com/RausellLab/NCBoost (22). The dataset is gene-centric, and no further processing was done.

#### Curation of TSG, OG, and NG training sets

The training set contained 242 high-confidence TSGs and 240 high-confidence OGs without overlapping from the v.87 CGC database on the COSMIC website, as well as 4058 NGs obtained as follows. The initial set of NGs was obtained from Davoli *et al.* (9). However, this initial set is likely to include mislabeled genes. To address this, those that overlap with the following gene lists (18 June 2020) were excluded from this initial NG set: (i) Candidate Cancer Gene Database (67), (ii) CancerMine (32), (iii) a cancer gene list compiled by Chiu *et al.* (68), (iv) the genes (OncoScore > 21.09) in the OncoScore database (69), and (v) allOnco Cancer Gene List (v3 February 2017; http://www.bushmanlab.org/links/genelists). The final training gene sets are available in data file S2.

#### **Candidate mutational features**

The candidate mutational features were previously defined by Davoli *et al.* (9) and Tokheim *et al.* (7). In addition to these features, other gene-centric features were also collected. Features were categorized into the following classes: Genomics, Mutation, Epigenetics, and Phenotype, and additional details regarding these features can be found in data file S1. The feature IDs mentioned below correspond to data file S1.

Features 1 to 20 were quantified on the basis of the combined mutation data using the script provided by Davoli *et al.* (9). Further information for these features can be found in their paper (9). For features 1, 5, and 6 in data file S1 that quantify the density of different categories of mutations within genes, only the coding sequence (CDS) length (per kilobase) of each gene is considered. For mutational features 8 to 15 and 28 that include ratios, a pseudocount estimated as the median of each feature across all genes was added, as described by Davoli *et al.* (9).

Features 11 to 15 rely on the functional effects of missense mutations, including high functional impact (HiFI) or low functional impact (LoFI) (9). The PolyPhen-2 Hum-Var prediction model was used to estimate the functional effects of missense mutations and to classify them as either HiFI or LoFI (9), based on the probability of functional damage as estimated by the PolyPhen-2 HumVar algorithm. Features based on HiFI and LoFI include the following: (i) benign mutations: silent and LoFI missense mutations; (ii) LoF mutations: nonsense and frameshift mutations; and (iii) HiFI missense mutations (damaging missense mutations). PolyPhen-2 scores (feature 16) were calculated by the PolyPhen-2 web server (http://genetics.bwh.harvard.edu/pph2/) (70). The missense MGAentropy scores (feature 33), which also measure the multispecies conservation of missense mutation sites, were also calculated by the Cancer-Related Analysis of Variants Toolkit (CRAVAT) tool (71).

Other mutation types include splicing/total mutations (feature 19) and inactivating fraction (feature 27). Splicing mutations are those that affect splicing sites; >95% of splicing mutations are in the first two positions at donor or acceptor sites. Inactivating mutations include indel frameshift, splice site, translation start site, and nonstop mutations. Features 21 to 29 that were introduced in Tokheim *et al.*'s paper were quantified based on our revised version of the script provided by Davoli *et al.* (9), given that these features can be quantified in a similar way to that for features 1 to 20. The lost start and stop fraction (feature 26) was defined as the fraction of the translation start site and nonstop mutations in total mutations. The recurrent missense fraction (feature 23) was defined by missense mutations occurring in more than one patient sample.

Features 42 to 46 are population genetics-based mutational features. For LoF constraints, three categories of tolerance to LoF mutations were defined by gnomAD: null (LoF mutations are fully tolerant), recessive (heterozygous LoF mutations are tolerant), and haploinsufficient (heterozygous LoF mutations are intolerant). The probability of the three types of mutations can also be obtained from the dataset (features 42 and 43) or be derived based on simple calculation (sum of the probability of three categories of intolerance equals 1). A pLI score was initially introduced to determine the likelihood that a given gene is intolerant of LoF mutations. The difference between LoF o/e and pLI is explained at https://blog. limbus-medtec.com/how-to-use-gnomad-v2-1-for-variant-filteringd7d2a7ee710a. For synonymous, missense, and LoF mutations (features 44 to 46), a signed Z score to describe the o/e was obtained from the gnomAD dataset. Higher Z scores indicate intolerance to variation or increased constraint, whereas lower Z scores indicate tolerance to variants.

#### Candidate epigenetic features

In addition to genetic data, epigenetic data have been shown to be associated with cancer driver genes. Here, we used the peak length and height to characterize histone modifications. We also used cancer-normal methylation difference to characterize gene promoter and gene-body methylation in cancer and normal samples. These potentially useful features (features 39 and 40 and 54 to 75) were previously used in epigenetics studies, but to what extent these features are useful in predicting cancer driver genes is not systematically evaluated. The histone modification BED files were processed based on our previously published procedures (19). Briefly, adjacent peaks were merged when peaks are within 3 kb by the merge command from bedtools (https://bedtools.readthedocs.io/). Peaks overlapping with the longest transcript of a gene with at least 50% of peak length were assigned to that gene by bedmap function in the BEDOPS tool (https://bedops.readthedocs.io/) with the following parameters: --max-element --echo --fraction-map 0.5 --delim '\t' --skip-unmapped. Features of "mean peak length" were calculated based on the merged peaks. For features of "height of peaks," the maximum signal values (seventh column in BED 6 + 4 narrow peak files used in ENCODE) were used. Promoter and gene-body cancer-normal methylation difference features (features 39 and 40) were defined by the mean methylation level in cancer samples (Beta Value column in the dataset) minus that in normal samples (Avg Beta Value Normal column in the dataset) based on COSMIC 450 K methylation data. 450K probes were mapped to genes according to genomic coordinates (hg19). The gene-body canyon cancer/normal methylation ratio feature (feature 41) was inspired from a previous

study (21). The ratio value was determined by the mean methylation level in cancer samples divided by that in normal samples in TCGA WGBS methylation data. To make "gene-body cancernormal methylation difference" (feature 40) and "gene-body canyon cancer/normal methylation ratio" (feature 41) available to all genes, genes without applicable feature values were imputed as 0. Genes were linked to gene-body canyons by BEDOPS with the same parameters as shown above. The difference between features 40 and 41 is that feature 41 is only available to genes with gene-body methylation canyons defined by a previous study using TCGA WGBS data (21), while feature 40 is available for all genes with 450K probes. We previously used TCGA WGBS data to define feature 41 because WGBS methylation data have a significantly higher resolution than 450K methylation data, while we were unable to identify large DNA methylation canyons using COSMIC 450K data. For feature 34, Repli-seq BAM datasets were quantified by the featureCounts program (http://subread.sourceforge.net/) to assign BAM reads to gene bodies. Read counts were normalized on the basis of the sequencing depth of the BAM files, and the normalized read numbers were used to calculate the S50 score (66). Early replication timing (feature 34) was quantified by the S50 score. All bam data are assigned to different cell cycle stages  $(G_1, S_1, S_2, S_3, \text{ and } S_4)$  for the S50 score calculation. This score, which determines the median replication timing (from 0 to 1), was calculated based on the algorithm proposed by a previous study (66). An S50 score that closes to 0 means early replication timing, whereas an S50 score that closes to 1 means late replication timing. Super enhancer percentage (feature 38) was calculated as the percentage of cell lines in which super enhancers are associated with any transcripts of genes.

#### Other candidate features

Feature 29 (log gene length) was defined as the log<sub>2</sub>-transformed length of the maximum transcript of a specific gene based on the ENSEMBL GTF annotation file. Feature 30 (log CDS length) was obtained from the COSMIC mutation files supplemented by the ENSEMBL GTF annotation file and then log<sub>2</sub>-transformed. Features 31 is CNA deletion percentage that was calculated based on column 17 (Mut type: gain or loss) in the original dataset (CNA amplification percentage can be calculated by 1 – CNA deletion percentage). The VEST scores (feature 35), which indicate missense pathogenicity for each mutation, were calculated by CRAVAT. Gene expression Z score (feature 36) was used to quantify the gene expression based on the "regulation" column in the original data. The exon conservation (phastCons) score (feature 32) that is based on the average phast-Cons score for maximum transcripts of genes was also calculated by CRAVAT. Feature 37 (increase of cell proliferation by CRISPR knockdown) was calculated on the basis of the cell proliferation scores in the CRISPR-screening data. A lower cell proliferation for a gene in a cell line means that the gene is more likely to be essential to the cell line. A score of 0 means nonessential, whereas a score of -1 means essential.

Features 47 to 53 are evolution-based features, including GDI (mutational damage that has accumulated in the general population), Primate dN/dS score (ratio between the number of nonsynonymous substitutions and the number of synonymous substitutions), RVIS percentile (high RVIS percentiles reflect genes highly tolerant to variation), ncRVIS, ncGERP, family member count (number of human paralogs for each gene), and the gene age (time of evolutionary origin based on the presence/absence of orthologs in vertebrates).

Genes with higher GERP scores are more constrained. ncRVIS is a measure of deviation from the genome-wide variants found in non-CDSs of genes (46). A negative ncRVIS score indicates less common noncoding variation than predicted. In ncRVIS and ncGERP, the non-coding regions were defined as the untranslated regions as well as nonexonic 250 bp upstream of TSSs.

#### Training of DORGE-TSG and DORGE-OG

The elastic net is a penalized regression method that can select a limited number of features that contribute to the response. Similar to the lasso, the elastic net selects features by shrinking some of the coefficients to be zero; the remaining features with nonzero coefficients are considered to have larger effects on the response and thus are selected and kept in the model. The main advantage of the elastic net over the lasso is that, in case of collinearity, the elastic net simultaneously selects a group of collinear features whereas the lasso tends to select only one feature from the group. (The simultaneous selection of collinear features is desired because, in the extreme situation where these collinear features are exactly identical, the regression method should assign equal coefficients to these features.) Therefore, we chose the elastic net over the lasso because we observed high collinearity among the original list of 75 features.

Specific to DORGE, we used LR with the elastic net penalty to train two binary classifiers for predicting TSGs and OGs, and these classifiers were referred to as DORGE-TSG and DORGE-OG. We used the R function glmnet from the R package glmnet (https://cran.r-project.org/web/packages/glmnet/index.html). The  $\lambda$  tuning parameter was selected by fivefold cross-validation using the function cv.glmnet from the same R package, while the  $\alpha$  parameter, which balances the lasso and ridge penalties, was set to the default value 0.5.

For every gene, DORGE-TSG predicted it with a probability of being a TSG, and this probability is defined as the gene's TSG score. The OG scores are defined similarly by DORGE-OG for all genes. Having two separate binary classifiers, one for detecting TSG and the other for detecting OG, DORGE is able to detect dual-functional genes.

The codes for training DORGE-TSG and DORGE-OG and obtaining predicted TSGs and OGs are available at https://github.com/biocq/DORGE. An online video that explains the code is available at www.youtube.com/watch?v=Pk8ZqoHK8zk.

#### **PRC** analyses

PRC analyses were performed using the R PRROC. The AUPRCs were calculated using TSG scores and OG scores by the pr.curve function in the package.

#### Thresholds on TSG scores and OG scores

We used the in-house code available in our DORGE GitHub repository to find the cutoffs on TSG scores and OG scores such that the population FPRs (type I errors; for TSG prediction, the FPR is the conditional probability of misclassifying an NG as a TSG) were controlled under 1%. The code was an implementation of the Neyman-Pearson classification umbrella algorithm (25).

## Gene sets and genomic and functional genomic datasets used for characterization and evaluation of DORGE-predicted novel TSGs and OGs

The gene lists and datasets that we used to evaluate our DORGE-predicted novel TSGs/OGs are as follows: (i) CGC gold-standard

gene list: The CGC is a widely used gold-standard list of cancerrelated genes. We used the CGC v.87 gene list as the testing gene set while excluding those in the CGC v.77 gene list to evaluate the performance of our prediction. (ii) ATAC-seq data: ATAC-seq data were taken from pan-cancer peak calls from data file S2 in Corces et al.'s paper (39). (iii) ERs: The ER gene list comes from a recent study focused on the characterization of ERs (40) and the EpiFactors database (72), after removing the genes that function only as TFs. (iv) Candidate TSGs identified by SB insertional mutagenesis. The inactivating pattern gene list was downloaded from the SBCDDB (42). This database contains cancer driver genes that were identified by SB insertional mutagenesis in tumor models. For the evaluation of DORGE-predicted novel TSGs, only genes with an inactivating pattern in the SBCDDB were kept, resulting in 1211 genes. (v) Survival data: Survival data were downloaded from the OncoLnc website (44). (vi) shRNA screening data. The gene-centric shRNA screening data were taken from the Achilles project (43). (vii) Evolutionary conservation data: For evolutionary conservation, we used phyloP scores that measure nonneutral substitution rates based on multispecies alignments. The phyloP data were downloaded from the UCSC (http://hgdownload.cse.ucsc.edu/goldenPath/ hg19/phyloP46way). We computed the average -log(phyloP) and phastCons score for each gene by averaging the base pair-level conservation values for every position in each gene. (viii) Phyletic age: We downloaded the precomputed phyletic age gene lists in human and measured enrichment of our predicted genes within the gene sets from different phyletic ages (i.e., Eukaryota, Metazoa, Chordata, and Mammalia) from the OGEE database (48). (ix) The BioGRID v3.5.183 data were downloaded from the website https://thebiogrid.org/. Biological network-related metrics can be calculated by the Cytoscape software (73). Additional information on the network metrics can be found in the Supplementary Text. (x) The PharmacoDB (74) gene-drug network data were downloaded from https://pharmacodb. pmgenomics.ca/. (xi) HKGs: We downloaded an HKG gene list from www.tau.ac.il/~elieis/HKG/, which includes 3804 HKGs. (xii) Essential genes: The essential and nonessential gene lists were also downloaded from the OGEE database. To shorten this list, we limited our essential gene set to those with >2 in entries of the OGEE database, resulting in 2340 definitive essential genes. Nonessential genes that overlap with essential genes were removed, resulting in 11,990 nonessential genes. (xiii) The drug response data were downloaded from Drug Gene Budger (54). Only significant drug-gene relationships (Q value < 0.05 and fold change > 2) were selected from the CREEDS data collections downloaded from the Drug Gene Budger database (54).

#### Gene set enrichment analysis

KEGG enrichment analyses were performed using Enrichr (75) for DORGE and DORGE variant predicted novel genes.

#### PPI network module analysis

For DORGE-predicted novel genes and CGC genes, PPI module analysis was performed by Metascape (76). Networks contain proteins that display physical interactions with at least one other protein in the list. For networks containing 3 to 500 proteins, the MCODE algorithm (77) was applied to identify densely connected network modules.

#### Statistical analysis

One-sided Wilcoxon rank-sum test was used when comparing different categories of genes. Gene enrichment analyses were performed in R, using one-sided Fisher's exact test (fisher.test function in R). *P* values of Spearman correlation were calculated by test for association/correlation between paired samples (cor.test function in R). Binomial test was used to test the enrichment of dual-functional genes in network hub genes.

#### **SUPPLEMENTARY MATERIALS**

Supplementary material for this article is available at http://advances.sciencemag.org/cgi/content/full/6/46/eaba6784/DC1

View/request a protocol for this paper from Bio-protocol.

#### **REFERENCES AND NOTES**

- 1. V. Labi, M. Erlacher, How cell death shapes cancer. Cell Death Dis. 6, e1675 (2015).
- B. Vogelstein, N. Papadopoulos, V. E. Velculescu, S. Zhou, L. A. Diaz Jr., K. W. Kinzler, Cancer genome landscapes. Science 339, 1546–1558 (2013).
- E. Y. H. P. Lee, W. J. Muller, Oncogenes and tumor suppressor genes. Cold Spring Harb. Perspect. Biol. 2, a003236 (2010).
- F. M. Behan, F. Iorio, G. Picco, E. Gonçalves, C. M. Beaver, G. Migliardi, R. Santos, Y. Rao, F. Sassi, M. Pinnelli, R. Ansari, S. Harper, D. A. Jackson, R. M. Rae, R. Pooley, P. Wilkinson, D. van der Meer, D. Dow, C. Buser-Doepner, A. Bertotti, L. Trusolino, E. A. Stronach, J. Saez-Rodriguez, K. Yusa, M. J. Garnett, Prioritization of cancer therapeutic targets using CRISPR-Cas9 screens. *Nature* 568, 511–516 (2019).
- S. A. Forbes, D. Beare, N. Bindal, S. Bamford, S. Ward, C. G. Cole, M. Jia, C. Kok, H. Boutselakis, T. De, Z. Sondka, L. Ponting, R. Stefancsik, B. Harsha, J. Tate, E. Dawson, S. Thompson, H. Jubb, P. J. Campbell, COSMIC: High-resolution cancer genetics using the catalogue of somatic mutations in cancer. *Curr. Protoc. Hum. Genet.* 91, 10.11.11–10.11.37 (2016).
- K. Tomczak, P. Czerwińska, M. Wiznerowicz, The Cancer Genome Atlas (TCGA): An immeasurable source of knowledge. *Contemp Oncol (Pozn)* 19, A68–A77 (2015).
- C. J. Tokheim, N. Papadopoulos, K. W. Kinzler, B. Vogelstein, R. Karchin, Evaluating the evaluation of cancer driver genes. *Proc. Natl. Acad. Sci. U.S.A.* 113, 14330–14335 (2016).
- M. H. Bailey, C. Tokheim, E. Porta-Pardo, S. Sengupta, D. Bertrand, A. Weerasinghe,
  A. Colaprico, M. C. Wendl, J. Kim, B. Reardon, P. K.-S. Ng, K. J. Jeong, S. Cao, Z. Wang,
  J. Gao, Q. Gao, F. Wang, E. M. Liu, L. Mularoni, C. Rubio-Perez, N. Nagarajan, I. Cortés-Ciriano,
  D. C. Zhou, W.-W. Liang, J. M. Hess, V. D. Yellapantula, D. Tamborero, A. Gonzalez-Perez,
  C. Suphavilai, J. Y. Ko, E. Khurana, P. J. Park, E. M. Van Allen, H. Liang; MC Working Group;
  Cancer Genome Atlas Research Network, M. S. Lawrence, A. Godzik, N. Lopez-Bigas, J. Stuart,
  D. Wheeler, G. Getz, K. Chen, A. J. Lazar, G. B. Mills, R. Karchin, L. Ding, Comprehensive
  characterization of cancer driver genes and mutations. Cell 173, 371–385.e18 (2018).
- T. Davoli, A. W. Xu, K. E. Mengwasser, L. M. Sack, J. C. Yoon, P. J. Park, S. J. Elledge, Cumulative haploinsufficiency and triplosensitivity drive aneuploidy patterns and shape the cancer genome. *Cell* 155, 948–962 (2013).
- M. Hofree, H. Carter, J. F. Kreisberg, S. Bandyopadhyay, P. S. Mischel, S. Friend, T. Ideker, Challenges in identifying cancer genes by analysis of exome sequencing data. *Nat. Commun.* 7, 12096 (2016).
- S. D. Bailey, K. Desai, K. J. Kron, P. Mazrooei, N. A. Sinnott-Armstrong, A. E. Treloar, M. Dowar, K. L. Thu, D. W. Cescon, J. Silvester, S. Y. C. Yang, X. Wu, R. C. Pezo, B. Haibe-Kains, T. W. Mak, P. L. Bedard, T. J. Pugh, R. C. Sallari, M. Lupien, Noncoding somatic and inherited single-nucleotide variants converge to promote *ESR1* expression in breast cancer. *Nat. Genet.* 48, 1260–1266 (2016).
- S. C. Mack, H. Witt, R. M. Piro, L. Gu, S. Zuyderduyn, A. M. Stütz, X. Wang, M. Gallo, L. Garzia, K. Zayne, X. Zhang, V. Ramaswamy, N. Jäger, D. T. W. Jones, M. Sill, T. J. Pugh, M. Ryzhova, K. M. Wani, D. J. H. Shih, R. Head, M. Remke, S. D. Bailey, T. Zichner, C. C. Faria, M. Barszczyk, S. Stark, H. Seker-Cin, S. Hutter, P. Johann, S. Bender, V. Hovestadt, T. Tzaridis, A. M. Dubuc, P. A. Northcott, J. Peacock, K. C. Bertrand, S. Agnihotri, F. M. G. Cavalli, I. Clarke, K. Nethery-Brokx, C. L. Creasy, S. K. Verma, J. Koster, X. Wu, Y. Yao, T. Milde, P. Sin-Chan, J. Zuccaro, L. Lau, S. Pereira, P. Castelo-Branco, M. Hirst, M. A. Marra, S. S. Roberts, D. Fults, L. Massimi, Y. J. Cho, T. Van Meter, W. Grajkowska, B. Lach, A. E. Kulozik, A. von Deimling, O. Witt, S. W. Scherer, X. Fan, K. M. Muraszko, M. Kool, S. L. Pomeroy, N. Gupta, J. Phillips, A. Huang, U. Tabori, C. Hawkins, D. Malkin, P. N. Kongkham, W. A. Weiss, N. Jabado, J. T. Rutka, E. Bouffet, J. O. Korbel, M. Lupien, K. D. Aldape, G. D. Bader, R. Eils, P. Lichter, P. B. Dirks, S. M. Pfister, A. Korshunov, M. D. Taylor, Epigenomic alterations define lethal CIMP-positive ependymomas of infancy. *Nature* 506, 445–450 (2014).
- S. M. Lauberth, T. Nakayama, X. Wu, A. L. Ferris, Z. Tang, S. H. Hughes, R. G. Roeder, H3K4me3 interactions with TAF3 regulate preinitiation complex assembly and selective gene activation. Cell 152, 1021–1036 (2013).
- 14. X. D. Zhao, X. Han, J. L. Chew, J. Liu, K. P. Chiu, A. Choo, Y. L. Orlov, W.-K. Sung, A. Shahab, V. A. Kuznetsov, G. Bourque, S. Oh, Y. Ruan, H.-H. Ng, C.-L. Wei, Whole-genome mapping

#### SCIENCE ADVANCES | RESEARCH ARTICLE

- of histone H3 Lys4 and 27 trimethylations reveals distinct genomic compartments in human embryonic stem cells. *Cell Stem Cell* 1, 286–298 (2007).
- M. J. Ziller, H. Gu, F. Müller, J. Donaghey, L. T.-Y. Tsai, O. Kohlbacher, P. L. De Jager, E. D. Rosen, D. A. Bennett, B. E. Bernstein, A. Gnirke, A. Meissner, Charting a dynamic DNA methylation landscape of the human genome. *Nature* 500, 477–481 (2013).
- Y. Bergman, H. Cedar, DNA methylation dynamics in health and disease. Nat. Struct. Mol. Biol. 20, 274–281 (2013).
- J. G. Herman, S. B. Baylin, Gene silencing in cancer in association with promoter hypermethylation. N. Engl. J. Med. 349, 2042–2054 (2003).
- R. Lister, M. Pelizzola, R. H. Dowen, R. D. Hawkins, G. Hon, J. Tonti-Filippini, J. R. Nery, L. Lee, Z. Ye, Q.-M. Ngo, L. Edsall, J. Antosiewicz-Bourget, R. Stewart, V. Ruotti, A. H. Millar, J. A. Thomson, B. Ren, J. R. Ecker, Human DNA methylomes at base resolution show widespread epigenomic differences. *Nature* 462, 315–322 (2009).
- K. Chen, Z. Chen, D. Wu, L. Zhang, X. Lin, J. Su, B. Rodriguez, Y. Xi, Z. Xia, X. Chen, X. Shi, Q. Wang, W. Li, Broad H3K4me3 is associated with increased transcription elongation and enhancer activity at tumor-suppressor genes. *Nat. Genet.* 47, 1149–1157 (2015).
- C. A. Davis, B. C. Hitz, C. A. Sloan, E. T. Chan, J. M. Davidson, I. Gabdank, J. A. Hilton, K. Jain, U. K. Baymuradov, A. K. Narayanan, K. C. Onate, K. Graham, S. R. Miyasato, T. R. Dreszer, J. S. Strattan, O. Jolanki, F. Y. Tanaka, J. M. Cherry, The Encyclopedia of DNA elements (ENCODE): Data portal update. *Nucleic Acids Res.* 46, D794–D801 (2018).
- J. Su, Y.-H. Huang, X. Cui, X. Wang, X. Zhang, Y. Lei, J. Xu, X. Lin, K. Chen, J. Lv, M. A. Goodell, W. Li, Homeobox oncogene activation by pan-cancer DNA hypermethylation. *Genome Biol.* 19, 108 (2018).
- B. Caron, Y. Luo, A. Rausell, NCBoost classifies pathogenic non-coding variants in Mendelian diseases through supervised learning on purifying selection signals in humans. *Genome Biol.* 20, 32 (2019).
- 23. A. Khan, X. Zhang, dbSUPER: A database of super-enhancers in mouse and human genome. *Nucleic Acids Res.* **44**, D164–D171 (2016).
- A. Tsherniak, F. Vazquez, P. G. Montgomery, B. A. Weir, G. Kryukov, G. S. Cowley, S. Gill, W. F. Harrington, S. Pantel, J. M. Krill-Burger, R. M. Meyers, L. Ali, A. Goodale, Y. Lee, G. Jiang, J. Hsiao, W. F. J. Gerath, S. Howell, E. Merkel, M. Ghandi, L. A. Garraway, D. E. Root, T. R. Golub, J. S. Boehm, W. C. Hahn, Defining a cancer dependency map. *Cell* 170, 564–576.e16 (2017).
- X. Tong, Y. Feng, J. J. Li, Neyman-Pearson classification algorithms and NP receiver operating characteristics. Sci. Adv. 4, eaao1659 (2018).
- D. Hnisz, B. J. Abraham, T. I. Lee, A. Lau, V. Saint-André, A. A. Sigova, H. A. Hoke, R. A. Young, Super-enhancers in the control of cell identity and disease. *Cell* 155, 934–947 (2013).
- M. G. Guenther, S. S. Levine, L. A. Boyer, R. Jaenisch, R. A. Young, A chromatin landmark and transcription initiation at most promoters in human cells. Cell 130, 77–88 (2007).
- C. R. Vakoc, M. M. Sachdeva, H. Wang, G. A. Blobel, Profile of histone lysine methylation across transcribed mammalian chromatin. Mol. Cell. Biol. 26, 9185–9195 (2006).
- L. A. Gates, J. Shi, A. D. Rohira, Q. Feng, B. Zhu, M. T. Bedford, C. A. Sagum, S. Y. Jung, J. Qin, M.-J. Tsai, S. Y. Tsai, W. Li, C. E. Foulds, B. W. O'Malley, Acetylation on histone H3 lysine 9 mediates a switch from transcription initiation to elongation. *J. Biol. Chem.* 292, 14456–14472 (2017).
- M. Zhao, P. Kim, R. Mitra, J. Zhao, Z. Zhao, TSGene 2.0: An updated literature-based knowledgebase for tumor suppressor genes. *Nucleic Acids Res.* 44, D1023–D1031 (2016).
- 31. Y. Liu, J. Sun, M. Zhao, ONGene: A literature-based database for human oncogenes. *J. Genet. Genomics* **44**, 119–121 (2017).
- J. Lever, E. Y. Zhao, J. Grewal, M. R. Jones, S. J. M. Jones, CancerMine: A literature-mined resource for drivers, oncogenes and tumor suppressors in cancer. *Nat. Methods* 16, 505–507 (2019).
- R. D. Kumar, A. C. Searleman, S. J. Swamidass, O. L. Griffith, R. Bose, Statistically identifying tumor suppressors and oncogenes from pan-cancer genome-sequencing data. *Bioinformatics* 31, 3561–3568 (2015).
- A. Gonzalez-Perez, N. Lopez-Bigas, Functional impact bias reveals cancer drivers. Nucleic Acids Res. 40, e169 (2012).
- N. D. Dees, Q. Zhang, C. Kandoth, M. C. Wendl, W. Schierding, D. C. Koboldt, T. B. Mooney, M. B. Callaway, D. Dooling, E. R. Mardis, R. K. Wilson, L. Ding, MuSiC: Identifying mutational significance in cancer genomes. *Genome Res.* 22, 1589–1598 (2012).
- F. Dietlein, D. Weghorn, A. Taylor-Weiner, A. Richters, B. Reardon, D. Liu, E. S. Lander, E. M. Van Allen, S. R. Sunyaev, Identification of cancer driver genes based on nucleotide context. *Nat. Genet.* 52, 208–218 (2020).
- D. Chakravarty, J. Gao, S. M. Phillips, R. Kundra, H. Zhang, J. Wang, J. E. Rudolph, R. Yaeger, T. Soumerai, M. H. Nissan, M. T. Chang, S. Chandarlapaty, T. A. Traina, P. K. Paik, A. L. Ho, F. M. Hantash, A. Grupe, S. S. Baxi, M. K. Callahan, A. Snyder, P. Chi, D. Danila, M. Gounder, J. J. Harding, M. D. Hellmann, G. Iyer, Y. Janjigian, T. Kaley, D. A. Levine, M. Lowery, A. Omuro, M. A. Postow, D. Rathkopf, A. N. Shoushtari, N. Shukla, M. Voss, E. Paraiso, A. Zehir, M. F. Berger, B. S. Taylor, L. B. Saltz, G. J. Riely, M. Ladanyi, D. M. Hyman, J. Baselga, P. Sabbatini, D. B. Solit, N. Schultz, OncoKB: A precision oncology knowledge base. JCO Precis. Oncol. 2017, (2017).

- S. S. Dhar, D. Zhao, T. Lin, B. Gu, K. Pal, S. J. Wu, H. Alam, J. Lv, K. Yun, V. Gopalakrishnan,
   E. R. Flores, P. A. Northcott, V. Rajaram, W. Li, A. Shilatifard, R. V. Sillitoe, K. Chen, M. G. Lee,
   MLL4 is required to maintain broad H3K4me3 peaks and super-enhancers at tumor suppressor genes. *Mol. Cell* 70, 825–841.e6 (2018).
- M. R. Corces, J. M. Granja, S. Shams, B. H. Louie, J. A. Seoane, W. Zhou, T. C. Silva, C. Groeneveld, C. K. Wong, S. W. Cho, A. T. Satpathy, M. R. Mumbach, K. A. Hoadley, A. G. Robertson, N. C. Sheffield, I. Felau, M. A. A. Castro, B. P. Berman, L. M. Staudt, J. C. Zenklusen, P. W. Laird, C. Curtis; The Cancer Genome Atlas Analysis Network, W. J. Greenleaf, H. Y. Chang, The chromatin accessibility landscape of primary human cancers. Science 362, eaav1898 (2018).
- L. Boukas, J. M. Havrilla, P. F. Hickey, A. R. Quinlan, H. T. Bjornsson, K. D. Hansen, Coexpression patterns define epigenetic regulators associated with neurological dysfunction. *Genome Res.* 29, 532–542 (2019).
- 41. F. Gnad, S. Doll, G. Manning, D. Arnott, Z. Zhang, Bioinformatics analysis of thousands of TCGA tumors to determine the involvement of epigenetic regulators in human cancer. BMC Genomics 16, (suppl 8), S5 (2015).
- J. Y. Newberg, K. M. Mann, M. B. Mann, N. A. Jenkins, N. G. Copeland, SBCDDB: Sleeping beauty cancer driver database for gene discovery in mouse models of human cancers. Nucleic Acids Res. 46, D1011–D1017 (2018).
- G. S. Cowley, B. A. Weir, F. Vazquez, P. Tamayo, J. A. Scott, S. Rusin, A. East-Seletsky, L. D. Ali, W. F. J. Gerath, S. E. Pantel, P. H. Lizotte, G. Jiang, J. Hsiao, A. Tsherniak, E. Dwinell, S. Aoyama, M. Okamoto, W. Harrington, E. Gelfand, T. M. Green, M. J. Tomko, S. Gopal, T. C. Wong, H. Li, S. Howell, N. Stransky, T. Liefeld, D. Jang, J. Bistline, B. H. Meyers, S. A. Armstrong, K. C. Anderson, K. Stegmaier, M. Reich, D. Pellman, J. S. Boehm, J. P. Mesirov, T. R. Golub, D. E. Root, W. C. Hahn, Parallel genome-scale loss of function screens in 216 cancer cell lines for the identification of context-specific genetic dependencies. Sci. Data 1, 140035 (2014).
- J. Anaya, OncoRank: A pan-cancer method of combining survival correlations and its application to mRNAs, miRNAs, and lncRNAs. PeerJ. Preprints 4, e2574v1 (2016).
- L. Liu, Y. Chang, T. Yang, D. P. Noren, B. Long, S. Kornblau, A. Qutub, J. Ye, Evolutioninformed modeling improves outcome prediction for cancers. *Evol. Appl.* 10, 68–76 (2017).
- S. Petrovski, A. B. Gussow, Q. Wang, M. Halvorsen, Y. Han, W. H. Weir, A. S. Allen,
   D. B. Goldstein, The intolerance of regulatory sequence to genetic variation predicts gene dosage sensitivity. *PLOS Genet.* 11, e1005492 (2015).
- K. W. Kinzler, B. Vogelstein, Cancer-susceptibility genes. Gatekeepers and caretakers. Nature 386, 761–763 (1997).
- W.-H. Chen, G. Lu, X. Chen, X.-M. Zhao, P. Bork, OGEE v2: An update of the online gene essentiality database with special focus on differentially essential genes in human cancer cell lines. *Nucleic Acids Res.* 45, D940–D944 (2017).
- A. A. Makashov, S. V. Malov, A. P. Kozlov, Oncogenes, tumor suppressor and differentiation genes represent the oldest human gene classes and evolve concurrently. Sci. Rep. 9, 16410 (2019).
- C. Cava, G. Bertoli, A. Colaprico, C. Olsen, G. Bontempi, I. Castiglioni, Integration of multiple networks and pathways identifies cancer driver genes in pan-cancer analysis. BMC Genomics 19, 25 (2018).
- H. Horn, M. S. Lawrence, C. R. Chouinard, Y. Shrestha, J. X. Hu, E. Worstell, E. Shea, N. Ilic,
   E. Kim, A. Kamburov, A. Kashani, W. C. Hahn, J. D. Campbell, J. S. Boehm, G. Getz, K. Lage,
   NetSig: Network-based discovery from cancer genomes. *Nat. Methods* 15, 61–66 (2018).
- D. Silverbush, S. Cristea, G. Yanovich-Arad, T. Geiger, N. Beerenwinkel, R. Sharan, Simultaneous integration of multi-omics data improves the identification of cancer driver modules. *Cell Syst* 8, 456–466.e5 (2019).
- E. Porta-Pardo, L. Garcia-Alonso, T. Hrabe, J. Dopazo, A. Godzik, A pan-cancer catalogue of cancer driver protein interaction interfaces. PLOS Comput. Biol. 11, e1004518 (2015).
- Z. Wang, E. He, K. Sani, K. M. Jagodnik, M. C. Silverstein, A. Ma'ayan, Drug Gene Budger (DGB): An application for ranking drugs to modulate a specific gene based on transcriptomic signatures. *Bioinformatics* 35, 1247–1248 (2019).
- A. S. Reddy, S. Zhang, Polypharmacology: Drug discovery for the future. Expert Rev. Clin. Pharmacol. 6, 41–47 (2013).
- The ICGC/TCGA Pan-Cancer Analysis of Whole Genomes Consortium, Pan-cancer analysis of whole genomes. Nature 578, 82–93 (2020).
- 57. M. S. Lawrence, P. Stojanov, P. Polak, G. V. Kryukov, K. Cibulskis, A. Sivachenko, S. L. Carter, C. Stewart, C. H. Mermel, S. A. Roberts, A. Kiezun, P. S. Hammerman, A. M. Kenna, Y. Drier, L. Zou, A. H. Ramos, T. J. Pugh, N. Stransky, E. Helman, J. Kim, C. Sougnez, L. Ambrogio, E. Nickerson, E. Shefler, M. L. Cortés, D. Auclair, G. Saksena, D. Voet, M. Noble, D. D. Cara, P. Lin, L. Lichtenstein, D. I. Heiman, T. Fennell, M. Imielinski, B. Hernandez, E. Hodis, S. Baca, A. M. Dulak, J. Lohr, D.-A. Landau, C. J. Wu, J. Melendez-Zajgla, A. Hidalgo-Miranda, A. Koren, S. A. McCarroll, J. Mora, R. S. Lee, B. Crompton, R. Onofrio, M. Parkin, W. Winckler, K. Ardlie, S. B. Gabriel, C. W. M. Roberts, J. A. Biegel, K. Stegmaier, A. J. Bass, L. A. Garraway, M. Meyerson, T. R. Golub, D. A. Gordenin, S. Sunyaev, E. S. Lander, G. Getz, Mutational heterogeneity in cancer and the search for new cancer-associated genes. *Nature* 499, 214–218 (2013).

#### SCIENCE ADVANCES | RESEARCH ARTICLE

- L. Mularoni, R. Sabarinathan, J. Deu-Pons, A. Gonzalez-Perez, N. López-Bigas, OncodriveFML: A general framework to identify coding and non-coding regions with cancer driver mutations. *Genome Biol.* 17, 128 (2016).
- D. Tamborero, A. Gonzalez-Perez, N. Lopez-Bigas, OncodriveCLUST: Exploiting the positional clustering of somatic mutations to identify cancer genes. *Bioinformatics* 29, 2238–2244 (2013).
- Z. Waks, O. Weissbrod, B. Carmeli, R. Norel, F. Utro, Y. Goldschmidt, Driver gene classification reveals a substantial overrepresentation of tumor suppressors among very large chromatin-regulating proteins. Sci. Rep. 6, 38988 (2016).
- J. Reimand, G. D. Bader, Systematic analysis of somatic mutations in phosphorylation signaling predicts novel cancer drivers. Mol. Syst. Biol. 9, 637 (2013).
- A. A. Stepanenko, Y. S. Vassetzky, V. M. Kavsan, Antagonistic functional duality of cancer genes. *Gene* 529, 199–207 (2013).
- 63. I. N. Smith, J. M. Briggs, Structural mutation analysis of PTEN and its genotype-phenotype correlations in endometriosis and cancer. *Proteins* **84**, 1625–1643 (2016).
- M. Lek, K. J. Karczewski, E. V. Minikel, K. E. Samocha, E. Banks, T. Fennell, A. H. O'Donnell-Luria, J. S. Ware, A. J. Hill, B. B. Cummings, T. Tukiainen, D. P. Birnbaum, J. A. Kosmicki, L. E. Duncan, K. Estrada, F. Zhao, J. Zou, E. Pierce-Hoffman, J. Berghout, D. N. Cooper, N. Deflaux, M. De Pristo, R. Do, J. Flannick, M. Fromer, L. Gauthier, J. Goldstein, N. Gupta, D. Howrigan, A. Kiezun, M. I. Kurki, A. L. Moonshine, P. Natarajan, L. Orozco, G. M. Peloso, R. Poplin, M. A. Rivas, V. Ruano-Rubio, S. A. Rose, D. M. Ruderfer, K. Shakir, P. D. Stenson, C. Stevens, B. P. Thomas, G. Tiao, M. T. Tusie-Luna, B. Weisburd, H.-H. Won, D. Yu, D. M. Altshuler, D. Ardissino, M. Boehnke, J. Danesh, S. Donnelly, R. Elosua, J. C. Florez, S. B. Gabriel, G. Getz, S. J. Glatt, C. M. Hultman, S. Kathiresan, M. Laakso, S. M. Carroll, M. I. McCarthy, D. M. Govern, R. M. Pherson, B. M. Neale, A. Palotie, S. M. Purcell, D. Saleheen, J. M. Scharf, P. F. Sklar, P. F. Sullivan, J. Tuomilehto, M. T. Tsuang, H. C. Watkins, J. G. Wilson, M. J. Daly, D. G. MacArthur; Exome Aggregation Consortium, Analysis of protein-coding genetic variation in 60,706 humans. *Nature* 536, 285-291 (2016).
- M. Jeong, D. Sun, M. Luo, Y. Huang, G. A. Challen, B. Rodriguez, X. Zhang, L. Chavez,
   H. Wang, R. Hannah, S.-B. Kim, L. Yang, M. Ko, R. Chen, B. Göttgens, J.-S. Lee, P. Gunaratne,
   L. A. Godley, G. J. Darlington, A. Rao, W. Li, M. A. Goodell, Large conserved domains of low
   DNA methylation maintained by Dnmt3a. *Nat. Genet.* 46, 17–23 (2014).
- C.-L. Chen, A. Rappailles, L. Duquenne, M. Huvet, G. Guilbaud, L. Farinelli, B. Audit, Y. d'Aubenton-Carafa, A. Arneodo, O. Hyrien, C. Thermes, Impact of replication timing on non-CpG and CpG substitution rates in mammalian genomes. *Genome Res.* 20, 447–457 (2010)
- K. L. Abbott, E. T. Nyre, J. Abrahante, Y.-Y. Ho, R. Isaksson Vogel, T. K. Starr, The Candidate Cancer Gene Database: A database of cancer driver genes from forward genetic screens in mice. *Nucleic Acids Res.* 43, D844–D848 (2015).
- H.-S. Chiu, S. Somvanshi, E. Patel, T.-W. Chen, V. P. Singh, B. Zorman, S. L. Patil, Y. Pan,
   S. S. Chatterjee; Cancer Genome Atlas Research Network, A. K. Sood, P. H. Gunaratne,
   P. Sumazin, Pan-cancer analysis of IncRNA regulation supports their targeting of cancer genes in each tumor context. *Cell Rep.* 23, 297–312.e12 (2018).
- R. Piazza, D. Ramazzotti, R. Spinelli, A. Pirola, L. De Sano, P. Ferrari, V. Magistroni,
   N. Cordani, N. Sharma, C. Gambacorti-Passerini, OncoScore: A novel, Internet-based tool to assess the oncogenic potential of genes. Sci. Rep. 7, 46290 (2017).
- I. Adzhubei, D. M. Jordan, S. R. Sunyaev, Predicting functional effect of human missense mutations using PolyPhen-2. Curr. Protoc. Hum. Genet. 7, 7.20.1–7.20.41 (2013).

- D. L. Masica, C. Douville, C. Tokheim, R. Bhattacharya, R. Kim, K. Moad, M. C. Ryan, R. Karchin, CRAVAT 4: Cancer-related analysis of variants toolkit. *Cancer Res.* 77, e35–e38 (2017)
- Y. A. Medvedeva, A. Lennartsson, R. Ehsani, I. V. Kulakovskiy, I. E. Vorontsov, P. Panahandeh, G. Khimulya, T. Kasukawa; The FANTOM Consortium, F. Drabløs, EpiFactors: A comprehensive database of human epigenetic factors and complexes. *Database* 2015, bav067 (2015).
- P. Shannon, A. Markiel, O. Ozier, N. S. Baliga, J. T. Wang, D. Ramage, N. Amin,
   B. Schwikowski, T. Ideker, Cytoscape: A software environment for integrated models of biomolecular interaction networks. *Genome Res.* 13, 2498–2504 (2003).
- P. Smirnov, V. Kofia, A. Maru, M. Freeman, C. Ho, N. El-Hachem, G.-A. Adam, W. Ba-alawi,
   Z. Safikhani, B. Haibe-Kains, PharmacoDB: An integrative database for mining *in vitro* anticancer drug screening studies. *Nucleic Acids Res.* 46, D994–D1002 (2018).
- M. V. Kuleshov, M. R. Jones, A. D. Rouillard, N. F. Fernandez, Q. Duan, Z. Wang, S. Koplev, S. L. Jenkins, K. M. Jagodnik, A. Lachmann, M. G. McDermott, C. D. Monteiro, G. W. Gundersen, A. Ma'ayan, Enrichr: A comprehensive gene set enrichment analysis web server 2016 update. *Nucleic Acids Res.* 44, W90–W97 (2016).
- Y. Zhou, B. Zhou, L. Pache, M. Chang, A. H. Khodabakhshi, O. Tanaseichuk, C. Benner, S. K. Chanda, Metascape provides a biologist-oriented resource for the analysis of systems-level datasets. *Nat. Commun.* 10, 1523 (2019).
- 77. G. D. Bader, C. W. V. Hogue, An automated method for finding molecular complexes in large protein interaction networks. *BMC Bioinformatics* **4**, 2 (2003).
- P. Chandrashekar, N. Ahmadinejad, J. Wang, A. Sekulic, J. B. Egan, Y. W. Asmann, S. Kumar, C. Maley, L. Liu, Somatic selection distinguishes oncogenes and tumor suppressor genes. *Bioinformatics* 36, 1712–1717 (2020).

Acknowledgments: We thank potential reviewers for their insightful suggestions and comments on this paper. Funding: This work was supported by grants from the U.S. NIH R01HG007538, R01CA193466, and R01CA228140 to W.L.; NIH R01GM120507, NSF grant DBI-1846216, Sloan Research Fellowship, Johnson & Johnson WiSTEM2D Award, and UCLA DGSOM W. M. Keck Foundation Junior Faculty Award to J.J.L.; and CPRIT RP160283-Baylor College of Medicine Comprehensive Cancer Training Program to F.P. Author contributions: Conceptualization: W.L., J.L., and J.J.L. Methodology: J.J.L. and J.L. Software: J.J.L. and J.L. Writing: J.J.L., J.L., Y.E.C., X.G., and W.L. Data analysis: J.L., J.J.L., J.S., F.P., and X.G. Supervision: W.L. and J.J.L. Competing interests: The authors declare that they have no competing interests. Data and materials availability: All data needed to evaluate the conclusions in the paper are present in the paper and/or the Supplementary Materials. Additional data related to this paper may be requested from the authors. The open source codes for DORGE are freely available at https://github.com/biocq/DORGE or upon request, along with the scripts to evaluate DORGE and reproduce the results of the paper.

Submitted 23 December 2019 Accepted 29 September 2020 Published 11 November 2020 10.1126/sciadv.aba6784

Citation: J. Lyu, J. J. Li, J. Su, F. Peng, Y. E. Chen, X. Ge, W. Li, DORGE: Discovery of Oncogenes and tumoR suppressor genes using Genetic and Epigenetic features. *Sci. Adv.* **6**, eaba6784 (2020).



### DORGE: Discovery of Oncogenes and tumoR suppressor genes using Genetic and Epigenetic features

Jie Lyu, Jingyi Jessica Li, Jianzhong Su, Fanglue Peng, Yiling Elaine Chen, Xinzhou Ge and Wei Li

Sci Adv **6** (46), eaba6784. DOI: 10.1126/sciadv.aba6784

ARTICLE TOOLS http://advances.sciencemag.org/content/6/46/eaba6784

SUPPLEMENTARY http://advances.sciencemag.org/content/suppl/2020/11/09/6.46.eaba6784.DC1 MATERIALS

**REFERENCES** This article cites 77 articles, 13 of which you can access for free

http://advances.sciencemag.org/content/6/46/eaba6784#BIBL

PERMISSIONS http://www.sciencemag.org/help/reprints-and-permissions

Use of this article is subject to the Terms of Service