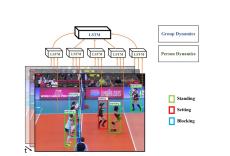
DeepFind: Sensor-driven Inference Acceleration for Continuous Deep **Mobile Vision Applications**

Chungkuk Yoo¹, Saiyma Sarmin², Inseok Hwang¹, Eric Rozner², Minsik Cho¹ ²University of Colorado Boulder ¹IBM

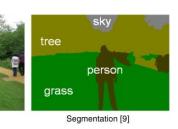
Problem and Goal

Continuous vision enables smart environments

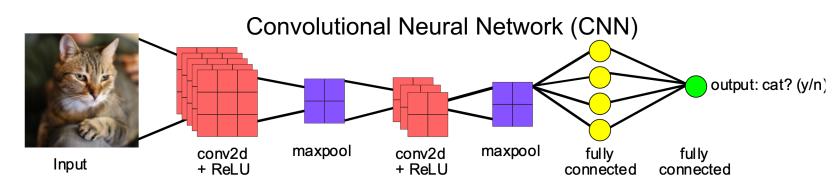




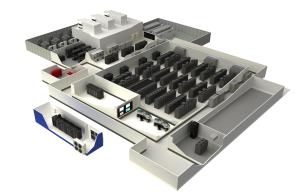




Deep learning CNNs obtain human-scale accuracy



Problem: CNN inference computationally expensive





- Privacy concerns
- Network cost





• Move computation to edge?

- Fewer resources than cloud (e.g., energy, computation)

Goal: enable deep learning vision to run continuously and efficiently on mobile and embedded devices.

Approach

Consecutive frames enable caching opportunities



Frame at t₀



Frame at t₁

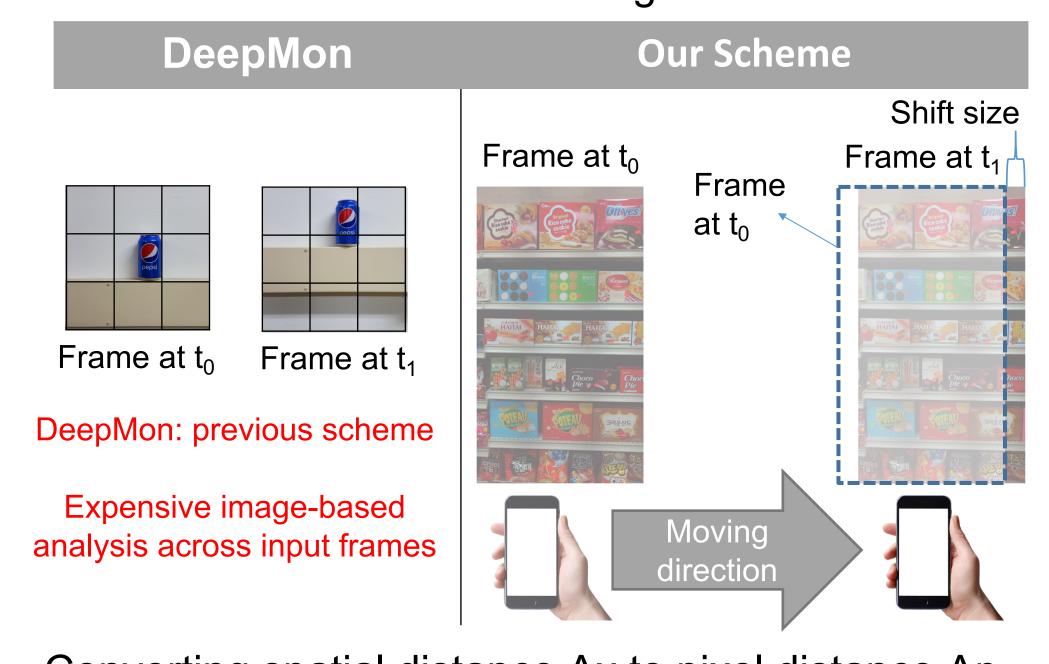


can be

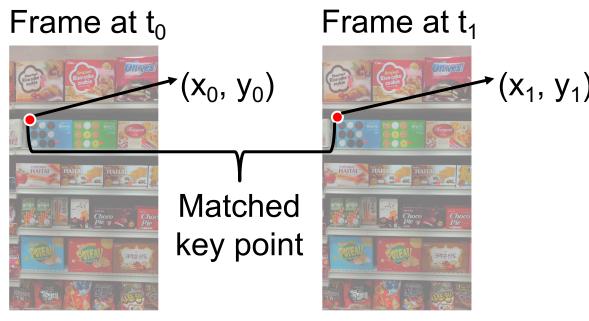
reused

Frame at t₀ Frame at t₁

• How to determine cacheable regions?



Converting spatial distance Δx to pixel distance Δp



Contributions

- Accelerate CNN on mobile and embedded devices
- A caching mechanism to reduce CNN inference time
- Exploits spatial/temporal similarities in CNN inputs

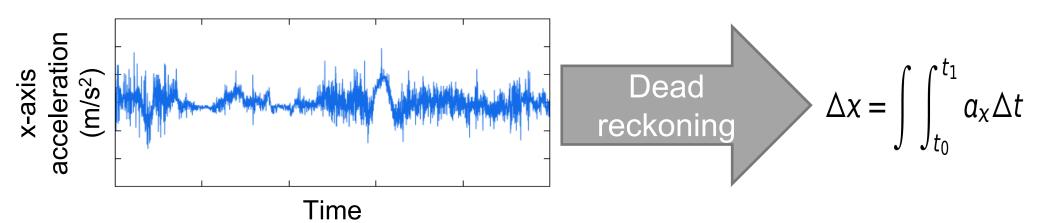








- Utilizes mobile sensors to determine similarities



Accelerator reading while moving

Phone displacement

Evaluation





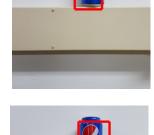




DeepFind



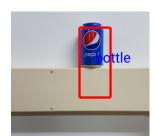




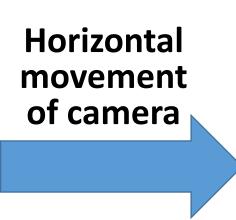
DeepMon

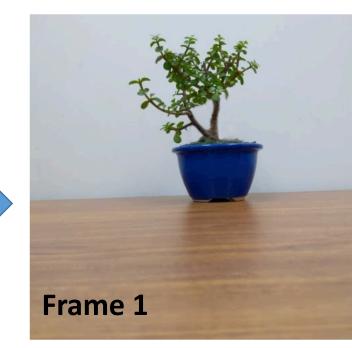


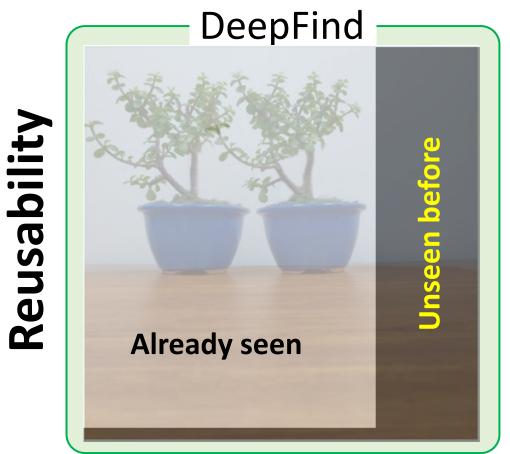


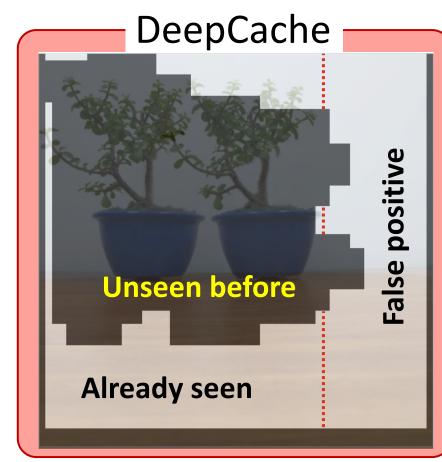












Time to determine cached region (per frame) DeepFind DeepMon DeepCache 11 - 30 ms0.42 ms 6.0 - 18 ms

Summary

- Continuous mobile vision important
 - Visual info provides context of users and environments
- Current deep learning algorithms are too expensive
- Edge devices have less power, energy than cloud Our work makes efficient continuous vision on mobile
- and embedded devices a reality
- Allows personalized intelligence to become truly pervasive