# Auto-reconfiguration for Latency Minimization in CPU-based DNN Serving

Ankit Bhardwaj MIT CSAIL

### Motivation

## Serving small DNN models on CPU:

- Consumes lower power
- Cost-effective
- Easy to scale across multiple machines

Increasing core counts and memory bandwidth make them suitable for serving small models.

Identifying Configuration is Challenging

<1, T, B> to <*i*, *t*, *b*> is challenging

- > Model specific parallelism
- Batch-specific gains
- > Heterogeneous hardware

Identifying a configuration that works for a specific model and a given batch size on a machine can be challenging.



Amar Phanishayee

Meta











Diminishing returns on adding more cores



Profile a few configurations for each model Predict the rest using a 2D-Knapsack solution

### Automated Reconfiguration



Deepak Narayanan **NVIDIA** 

Filled with DP Algo







### Conclusion

Pure intra-op parallelism on CPUs shows diminishing returns.

Multi-instance execution improves

Packrat automatically finds optimal configurations and dynamically

adjusts the serving system.

> Achieves up to **3.2x** latency speedup.