# Serving and Optimizing Machine Learning Workflows on Heterogeneous Infrastructures

Yongji Wu, Matthew Lentz, Danyang Zhuo, Yao Lu



# 

# ML Workflows are Widely Deployed



## **Multi-Camera Vehicle Tracking**







#### **Visual Question Answering**



# Growth of Heterogeneous Infrastructure

- Inputs derived from sensors at the edge
  - Cameras, IoT devices, etc
  - Some even have support for on-board compute



**Boulder Al DNN** Cam

- **Connected to more powerful compute over network** 
  - Edge-local hubs and cloud computing resources



# Many Choices for Serving Workflows

## Worker Assignment (WA)



## Model Selection (MS)



Different architectures Model compression techniques

# JellyBean: Optimizing Execution Plans

Goal

- - Each decision is affected by the
  - Should yield more efficient exect
    - However, search space becomes

Prior work only leveraged WA —

## Minimize serving cost subject to throughput and accuracy constraints

## **Insight:** Jointly optimize worker assignment and model selection

other	System	Parallelism	Q MS	O WA
ution plans	PyTorch [52]	Data	×	×
	TF [12]	Data	×	×
s much larger!	Spark [65]	Data	×	×
	Clipper [21]	Data	×	×
	Ray [47]	Data, Model	×	×
	Optasia [42]	Data, Op	×	$\checkmark$
	Pathways [14]	Data, Model	×	×
	Llama [58]	Data, Op	×	$\checkmark$
	Scrooge [27]	Data, Op	×	$\checkmark$
	JellyBean (Ours)	Data, Op	$\checkmark$	$\checkmark$

# JellyBean: System Overview







- Profile models across workers to understand impact of model selection
- Use profiles (and constraints) to compute an optimized execution plan



# Overview of User Inputs



Infrastructure

Properties of the available compute and network resources: Type (e.g., CPU, V100 GPU), Count, Bandwidth, Cost



Targets for minimum accuracy and supported throughput

A directed acyclic graph (DAG) with nodes as inputs or operators (ML / relational)

# User Inputs: Model Choices



## Many model variants with different accuracy-efficiency tradeoffs

# **User Inputs: Infrastructure Tiers**

Heterogeneous (Across Tiers)



#### Homogeneous (Within Tiers)





Compute accuracy response for model based on direct upstream models 2 We assume output accuracy is monotonically increasing w.r.t. input accuracy



Acc = ?

# **Query Optimizer: Objective Function**



 $s.t. acc \ge A, t_{v_{out}}^{a(v_{out})} \ge T$ 

Assumptions

"Communication among workers within a tier is free"  $\rightarrow C_B$  set to 0 "Communication is one way from lower to upper tiers"  $\rightarrow C_B$  set to  $\infty$ 

#### **Communication Cost**

Based on bandwidth cost between workers **C**<sub>B</sub> and the consumed bandwidth **R** for output of one model as input to the other



# Query Optimizer: Approaches to MS and WA

## **Model Selection**

- Go in reverse topological order
  - Final accuracy must satisfy the user's minimum constraint
  - Keep top selections in the beam based on simplified cost model (assume most powerful worker)

## Worker Assignment

- Go in topological order
  - Reasonable approximation given realistic workflows (and one-way assumption)
- Greedy choice of lowest-cost workers (preferring lower tiers)
  - Keep top assignments in the beam based on accurate cost model (knowledge of worker)



# **Query Processor**



(Naiad)

## **Our Extensions for JellyBean:**

- Support for operator-level parallelism instead of data-level parallelism
- Support heterogeneous runtimes 2
- Support relational operators (e.g., filter, join) 3

## Timely Dataflow provides a low-overhead dataflow abstraction

# **Evaluation: Inputs**

## Workflows

## **NVIDIA AI City Challenge (AICity)**



## Visual Query Answering (VQA)



## Infrastructures

#### 5/9/15/30 workers (CPU & GPU) divided into multiple tiers



# Example: JellyBean Execution Plan for VQA

## "Medium" Infrastructure:



Constraints: Throughput  $\ge 20$  rps Accuracy  $\ge 0.56$ 



# **Evaluation: Total Serving Cost**



## **Baselines:** Note: All use most accurate models FF = First Fit (Cheapest worker on any tier) **BF = Best Fit** (Cheapest worker on current tier) **PTc = PyTorch on Cloud**

(One cloud GPU worker)

**SPc = Spark with PyTorch** (All cloud GPU workers)

## Ideal:

**LB** = Lower Bound (Brute force over all plans)

**PTe = PyTorch on Edge** (One edge GPU worker)



# **Evaluation: Serving Cost vs Throughput**



JellyBean (JB) can keep up with increasing throughputs (and is near optimal)

# Summary

## JellyBean: Automatically optimizes deployments for ML workflow serving

Jointly leverages worker assignment and model selection

Optimized plans can significantly reduce the total serving cost

![](_page_17_Picture_4.jpeg)

https://github.com/libertyeagle/JellyBean