

# CS 744: NEXUS

Shivaram Venkataraman

Fall 2022

# ADMINISTRIVIA

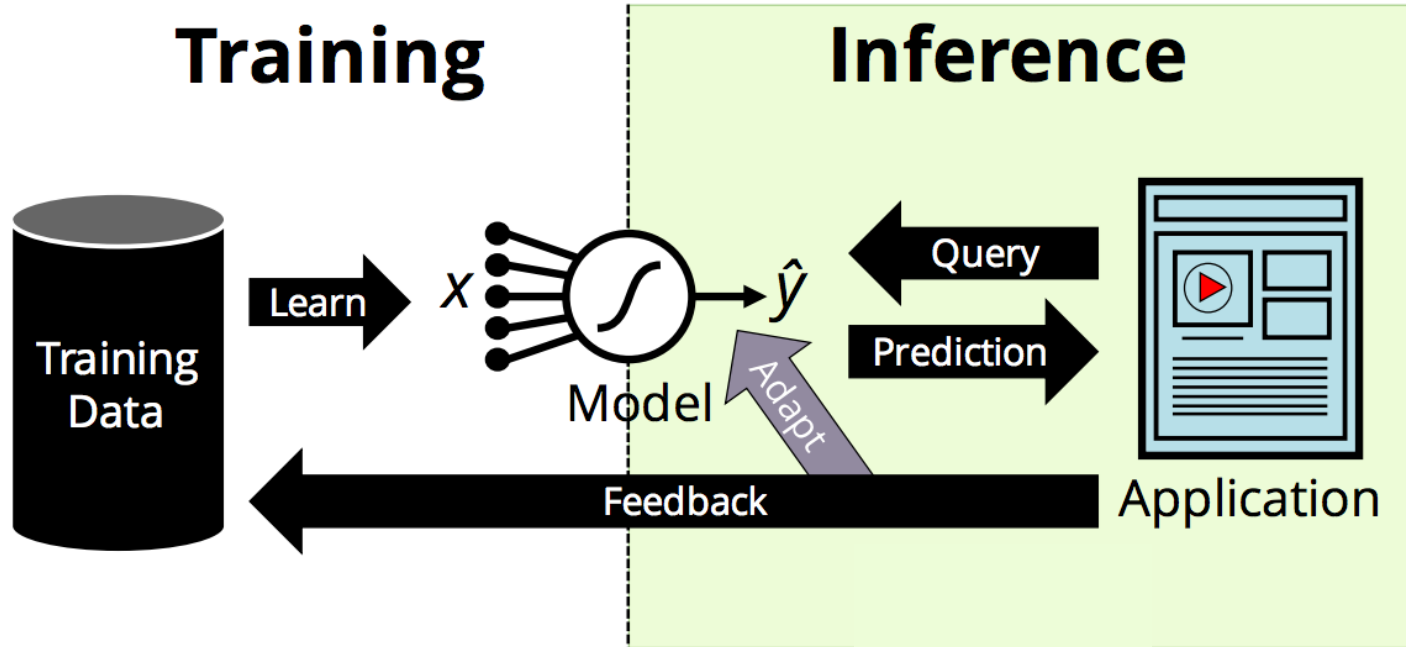
## Course Project Proposals

- Due Oct 26!
- See Piazza for template

## Midterm details

- Oct 27<sup>th</sup>: Includes papers from Datacenter as a Computer to Nexus
- Open book, open notes
- Held in class time 9.30-10.45am Central Time

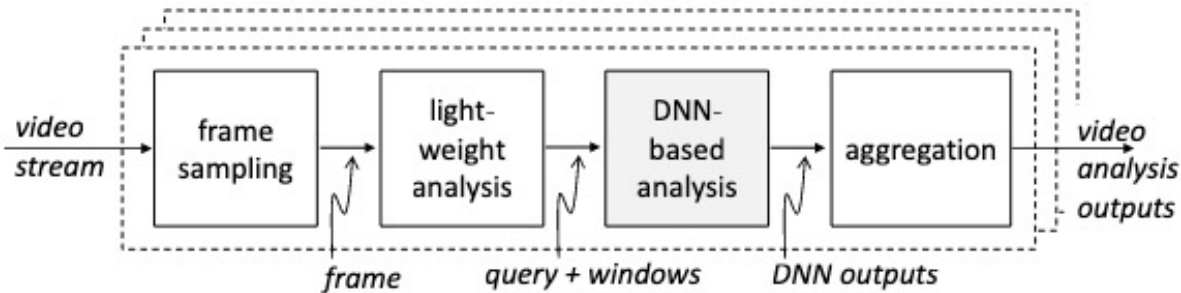
# MACHINE LEARNING: INFERENCE



# EXAMPLE APPLICATION

## Video analysis service

- Thousands of streams, thousands of tenants
- Each stream is processed by a DNN-based “query”
- Latency SLOs (10s to 100s of ms)



# SCHEDULING GOAL: HIGH GPU UTILIZATION

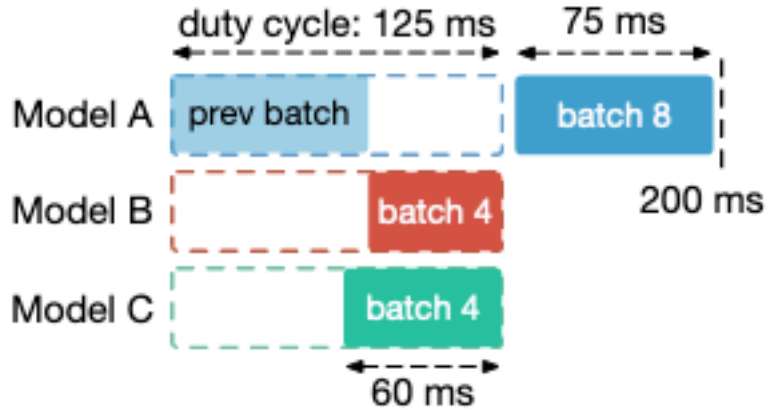
Placement

Batching

$$\text{batch\_lat}(b) = \alpha b + \beta,$$

# SCHEDULING BATCHED EXECUTION

Target tputs A: 64, B: 32, C: 32 req/sec. SLO: 250ms



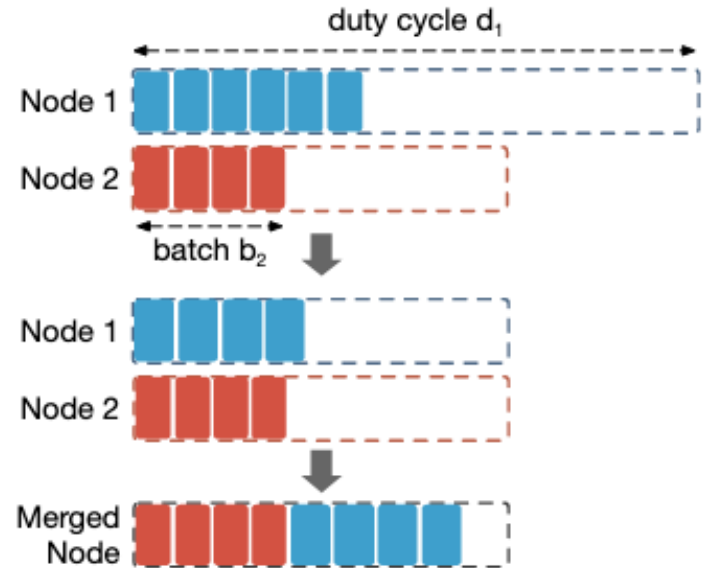
Model A			Model B			Model C		
Batch	Lat	Req/s	Batch	Lat	Req/s	Batch	Lat	Req/s
4	50	80	4	50	80	4	60	66.7
8	75	107	8	90	89	8	95	84
16	100	160	16	125	128	16	125	128

# BATCH-AWARE SCHEDULING

Inputs: Request rate, SLO for each model, Profiles at batch size

Approach: Allocate “full” GPUs based on load. Handle residuals

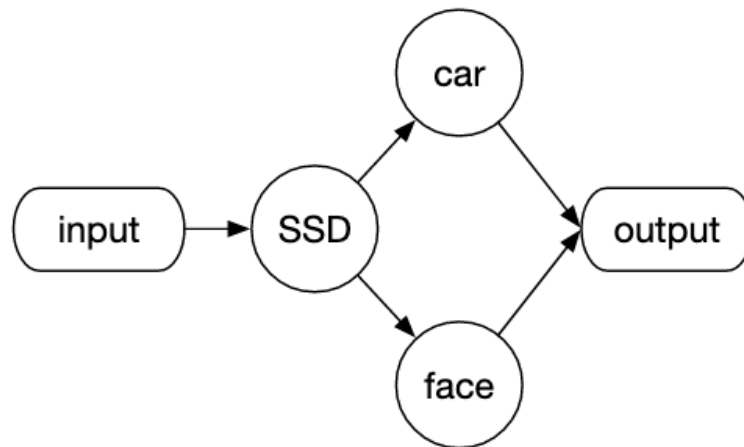
Greedy Approximation



# HANDLING COMPLEX QUERIES

Challenge:

How do we set latency SLOs for complex queries?





# SCHEDULING COMPLEX QUERIES

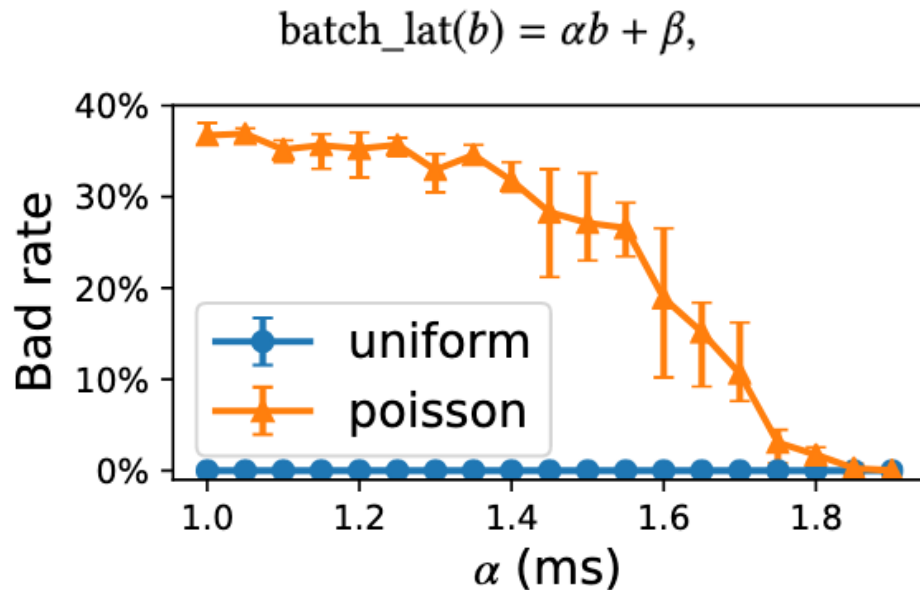
Query Analysis to determine latency SLO splits

Inputs: Models with request rate  $R_i$  latency SLO  $L$

$$\begin{aligned} & \underset{\{b_v\}}{\text{minimize}} && \sum_v R_v l_v(b_v) / b_v \\ & \text{subject to} && \sum_{u: M_{\text{root}} \rightsquigarrow M_v} l_u(b_u) \leq L \quad \forall v \in \text{leaf} \end{aligned}$$

# ADAPTIVE BATCHING

Clipper: Adapt the batch size based on the oldest request in the queue

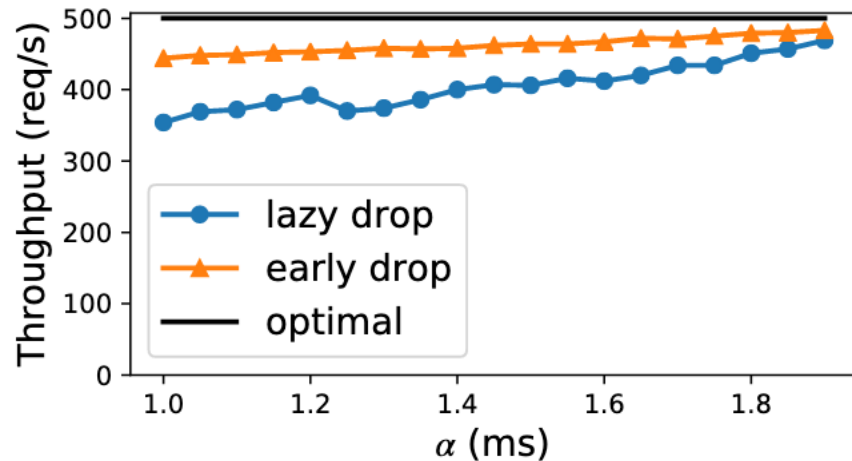


# BATCH-AWARE DISPATCH

Early-dropping scheme

1. Scans queue using sliding window of batch size

2. Stop at the first request with that can execute *entire window*



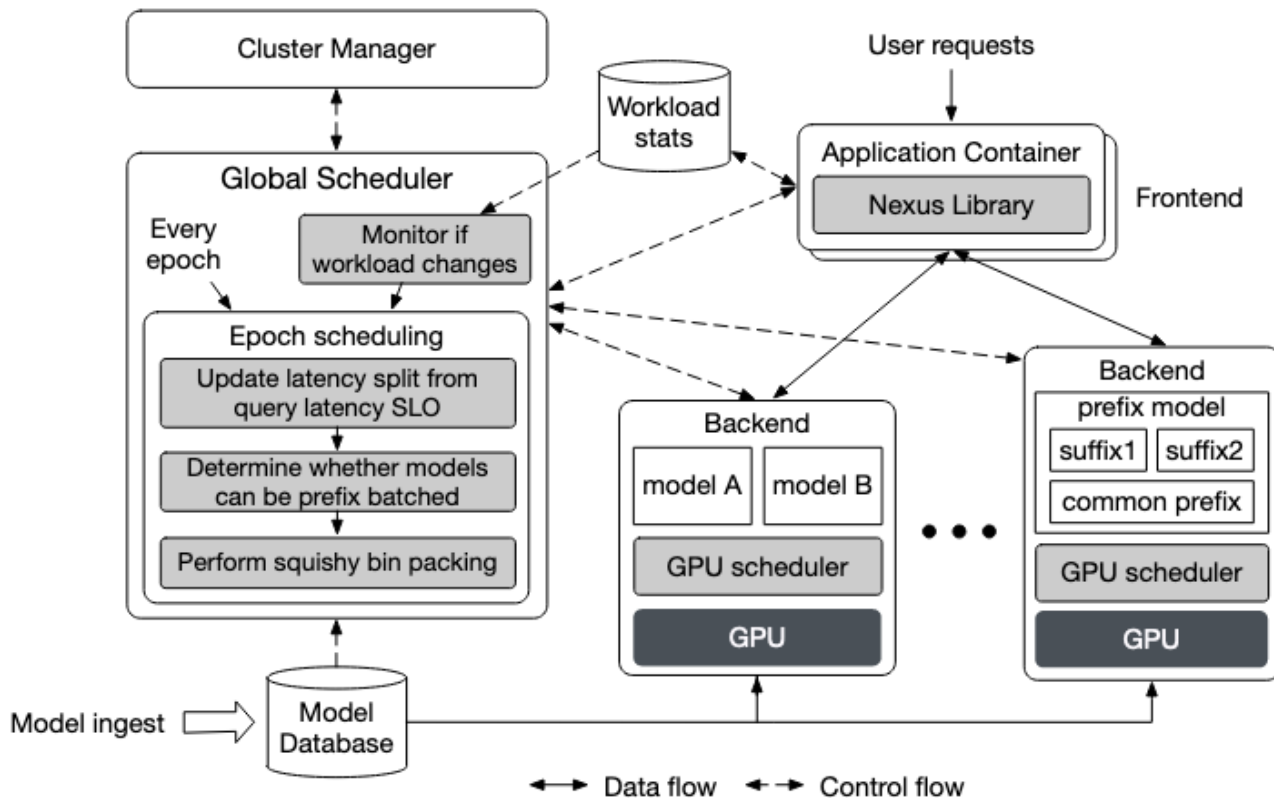
# OTHER FEATURES

Prefix Batching

GPU Multiplexing

Overlapping CPU and GPU computation

# NEXUS ARCHITECTURE



# SUMMARY

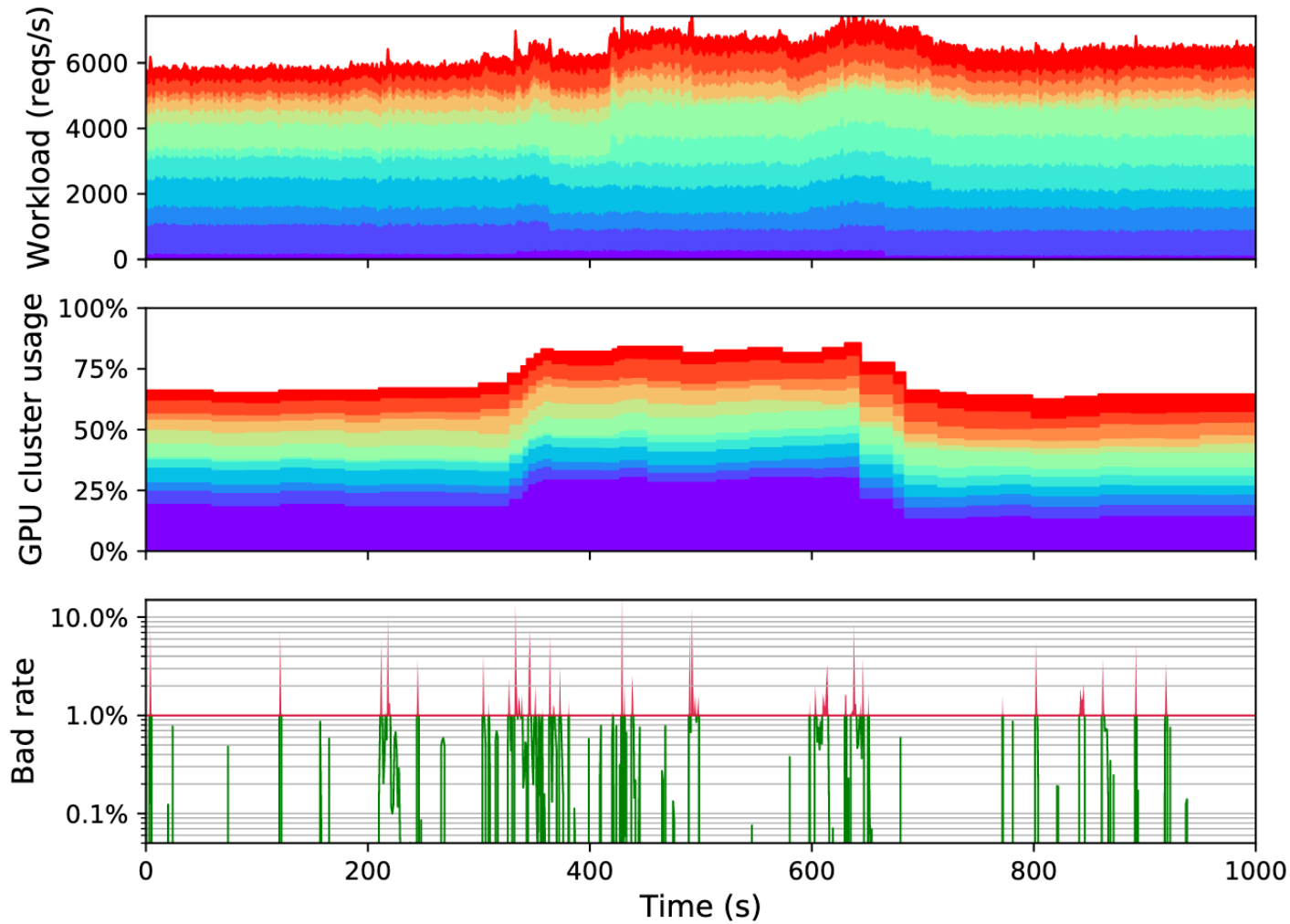
- ML Inference goals: latency SLO, GPU utilization
- Nexus: Handle multiple tenants, multiple DNNs
- Schedule using squishy bin packing
- Breakdown SLO for complex queries, adaptive batching

# DISCUSSION

<https://forms.gle/PtEaiF4casfZm2JY6>

Consider a scenario where you have a model that takes variable amount of time depending on the input. For example if a frame contains 100 cars it takes 250ms to process but if the frame has 1 car then it finishes in 10ms. What could be one shortcoming in using Nexus to schedule this model?





Next class: SQL

Coming soon

Project Introductions

Midterm I