Systems support for ML: It's all about the choices

Neeraja J. Yadwadkar UT Austin, ECE Nov 8th, 2022

With collaborators from Stanford and UT Austin

What is ML?





A dumb algorithm with lots and lots of data beats a clever one with a modest amount of it!*

* A Few Useful Things to Know About Machine Learning, CACM'12

Why is ML research flourishing?

- Availability of data
- Availability of compute
- Advances in algorithms and models

Abstractions and Interfaces!





Facebook's pipeline



Applied Machine Learning at Facebook: A Datacenter Infrastructure Perspective, HPCA'18

Google's pipeline

Integrated Frontend for Job Management, Monitoring, Debugging, Data/Model/Evaluation Visualization



Design Complexity

Machine Learning

- Model
- Feature engineering
- Training
- Bias
- Overfitting
- Generalizability
- Accuracy



Systems

- Scheduling
- Resource allocation
- Locality
- Fault tolerance
- Power efficiency
- Reliability/Availability
- Security

Systems for ML research brings these two sets of complexities together



Outline

➤ What is ML?

- > ML Workflow
- Systems for ML: Design Complexity
 - Managed Inference Serving by INFaaS (Deep dive)
- Inference Pipelines using Llama
- Sparsely Activated Massive Models





Same model, multiple *Model-Variants*!

Same model, multiple Model-Variants!

Compiler optimizations

TVM, TensorRT



Same model, multiple Model-Variants!

Compiler optimizations

TVM, TensorRT

Different precisions

INT8, FP16, FP32

Hyperparameter optimizations

Batch size

Heterogeneous Hardware

CPUs, GPUs, ..., *PUs

Trade-offs for inference due to heterogeneous hardware



Observations:

- CPUs: increased inference latency with batch-size
- **GPUs**:
 - Higher loading latencies
 - Perform significantly better on high batch sizes

Trade-offs for inference due to heterogeneous hardware



Same model, multiple *Model-Variants*!



We can compile the same model to 10s (100s) of versions



Default user choices



• Overprovision:

- Use dedicated resources
- Keep the models "always on"
- Replicate a query across multiple models

Model-Variants: Challenge, but an untapped Opportunity!

Challenge	Opportunity	
Hardware x Optimizers x Precisions x	A large trade-off space of Latency, Throughput,	Our Proposal: Model-Iess Inference Serving
Hyperparameters → a large search	Accuracy, Resources required, and Cost	
space		







Key: Automatically and efficiently selecting and scaling model-variants

INFaaS: A Model-less Inference Serving System

- > No models to generate and manage for users
- > Automatic selection of the right model-variant for each query
- > Autoscaling to respond to the changes in query load

INFaaS' Model-less API

INFaaS' Model Registration API

 INFaaS' Query Submission API

query(input.jpg, detectFaceApp, latency=200ms, accuracy=70%)

INFaaS: Architecture Overview



INFaaS: Model Registration Workflow



INFaaS: Query Execution Workflow



INFaaS: Query Execution Workflow



Selecting a Variant



INFaaS: Autoscaling Workflow



INFaaS' Model-Autoscaler



Question: What combination of variants (types and number) is required to support the changed load?

Outline

- ➤ What is ML?
- > ML Workflow
- Systems for ML: Design Complexity
- Managed Inference Serving by INFaaS (Deep dive)
 - Inference Pipelines using Llama
- Sparsely Activated Massive Models

Inference pipelines

- Real-world applications issue pipelines that include inference tasks
 - Example: Video pipelines are ubiquitous with various cost-perf targets





"Add a vintage filter to "Identify cars and faces the video" from the traffic feed"



Users configure operation knobs to best meet targets

<Hardware resources, batch size, resolution, ...>



Challenges:

- Large configuration space
- Input-dependent execution flow
- Exhaustive profiling is expensive



Llama



Fine-grained SLOs (E2E SLO to per-operation SLOs)

Fine-grained resource allocations (per-video, per-frame)

Outline

- ➤ What is ML?
- > ML Workflow
- Systems for ML: Design Complexity
- Managed Inference Serving by INFaaS (Deep dive)
- Inference Pipelines using Llama
 - Sparsely Activated Massive Models

Massive Neural Networks

- The capacity of a neural network to absorb information is limited by its number of parameters.
 - More capacity \rightarrow more accuracy
 - But more the parameters ightarrow more computation
- So sparse models have been proposed
 - Decreased accuracy
- Conditionally sparse models: parts of the network are activated per example
 - Improved model capacity without increasing number of parameters



*Outrageously large neural networks: The sparsely-gated Mixture-of-Experts layer, ICLR'17

Research questions

- Resource allocation for inference requests?
- Inference in resource-constrained edge settings

Systems support for ML: It's all about the choices

 ML is a new workload that imposes various new trade-offs and choices for systems

- Choices that matter
- Choices that are hard to make
- Who, better than systems folks, should navigate these choices?

Neeraja J.Yadwadkar neeraja@austin.utexas.edu