



## Hardware Opportunities in the Machine Learning Lifecycle

Joseph E. Gonzalez  
Co-director of the RISE Lab  
[jgonzal@cs.berkeley.edu](mailto:jgonzal@cs.berkeley.edu)

1

Get the latest slides  
and links to literature




<https://tinyurl.com/isscc-lifecycle>

2

### About Me

- Co-director of the RISE Lab
- Co-founder of Turi Inc.
- Member of the Apache Spark PMC
- Research
  - Artificial Intelligence
  - Data Science
  - Distributed Data Systems
  - Graph Processing Systems
- I don't study processor **architecture**
  - But I probably should ...



3

### Outline

- History and the **Co-evolution** of Hardware and AI
  - The **Feedback Cycle** driving the 3<sup>rd</sup> wave of AI
- Machine Learning is **not a single workload**
  - Stages of the **Machine Learning Lifecycle**
- Security** and machine learning

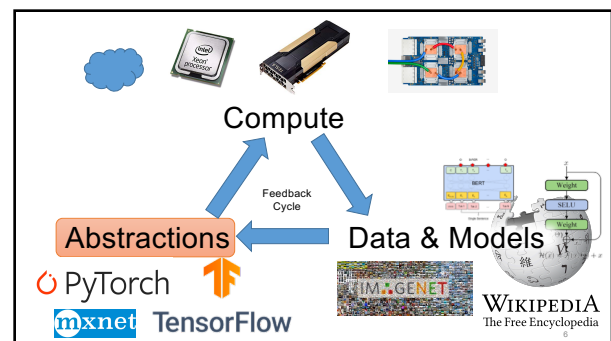
*Along the way, I will talk about some of the research in my group addressing interesting aspects of the lifecycle.*

4

### Hardware and the History of AI

- 1950 to 1974: Birth of AI**
  - 1951 Marvin Minsky builds first **neural network hardware** (SNARC)
- 1974 to 1980: First AI Winter**
  - Limited **processing power** and **data**
- 1980 to 1987: Second Wave of AI**
  - XCON (AI for **hardware configuration**) for DEC → boom in AI hardware companies
- 1987 to 1993: Second AI Winter**
  - Brittle** AI and the collapse of the **AI Hardware Market**
- 1993 to 2011: AI → Machine Learning**
  - Confluence of ideas + **Compute** + **Big Data** → AI starts to really work
- 2011 to 2019: Third Wave (Deep Learning)**
  - Compute** + **data** + **abstractions** → **feedback cycle**

5



## Abstractions are Enabling Innovation

- Much of machine learning before 2010
  - Research focused on machine learning **algorithms**
  - Programs written using **high-level imperative languages**
    - Matlab/R/C++/Java
  - **Big abstractions**: linear algebra, map-reduce, graph systems
- Today:
  - Research focused on **model design**
  - **Models** written in high-level **DSLs**
    - TensorFlow/Pytorch
  - Big abstractions: **tensor operations**, **loss minimization**, linear algebra, ...
- Models written in **TensorFlow** can now run on **hardware** that didn't exist when the models were created.

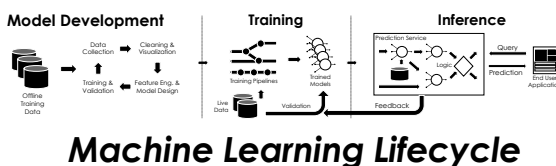
8

How do we make **hardware** for  
**Machine learning**?

**Machine learning**  
is not a single application.

9

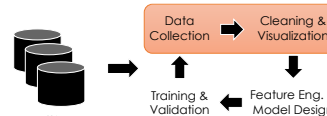
*Machine learning is  
multiple applications with different requirements.*



**Machine Learning Lifecycle**

10

## Model Development



Identifying potential  
sources of data

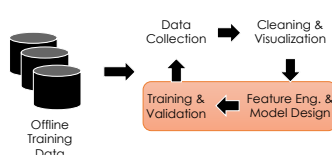
Joining data from  
multiple sources

Addressing **missing  
values and outliers**

**Plotting** trends to  
identify **anomalies**

11

## Model Development



Building informative  
**features functions**

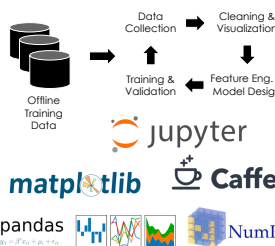
Designing new **model  
architectures**

**Tuning** training algos.

**Validating** prediction  
accuracy

12

## Model Development Frameworks



13

## Model Development → Hardware

- Need to test **multiple designs and hyperparameters** quickly
  - May be better to run many parallel experiments than one experiment faster
- **Debug heavy** → sources of error → data, hyperparams., & model
  - **System should not be a source of error**
  - **Avoid cutting corners** (e.g., quantization, async) for increased performance
    - Unless you can make a case for stable convergence ...
- **Data preparation** is often a bottleneck
  - Opportunity for **data tooling**
  - Accelerate data transformation and augmentation
- **Emerging Trends**
  - **Attention Models** and **Graph Neural Networks**: reduced locality, sparsity
  - **Dynamic Networks**: gating, cascades, mixtures, ...
  - Increased emphasis on **DNN features** and **Fine-Tuning**
    - Reuse of common architectures and weights

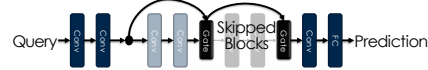
14

## Dynamic Networks for **fast** and **accurate** inference

**IDK Cascades:** Using the fastest model possible [UAI'18]

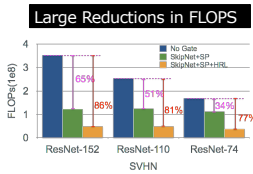
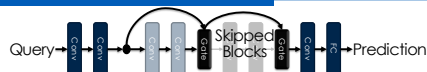


**SkipNet:** dynamic execution within a model [ECCV'18]



15

## SkipNet: dynamic execution within a model [ECCV'18]

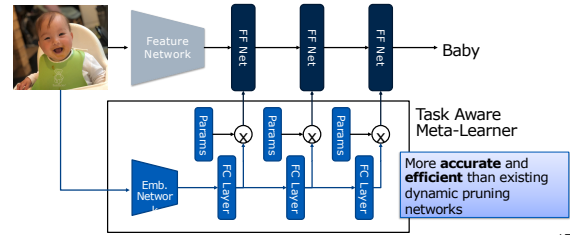


**Skip more layers on clear images**



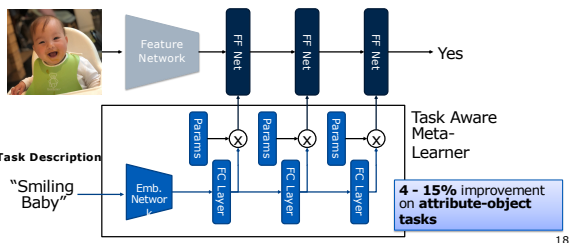
16

## Task Aware Feature Embeddings [CVPR'19]



17

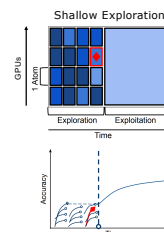
## Task Aware Feature Embeddings [CVPR'19]



18

## HyperSched [SOCC'19]

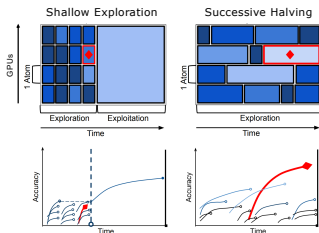
**Dynamically** allocating **parallel resources** to **parallel experiments**.



19

## HyperSched [SOCC'19]

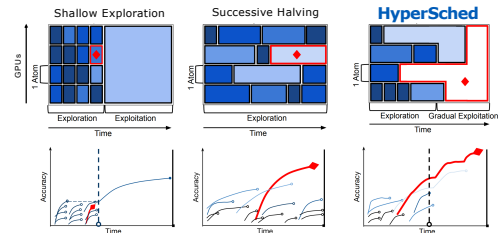
Dynamically allocating parallel resources to parallel experiments.



20

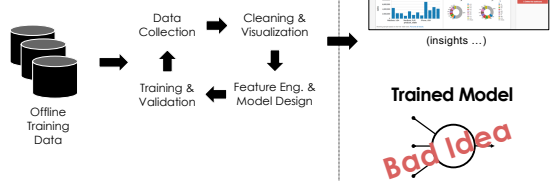
## HyperSched [SOCC'19]

Dynamically allocating parallel resources to parallel experiments.



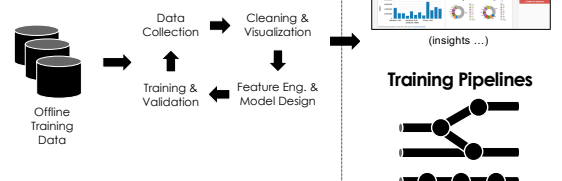
21

## What is the output of Model Development



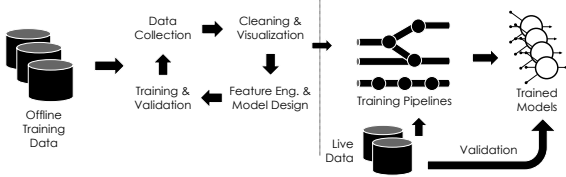
22

## What is the output of Model Development



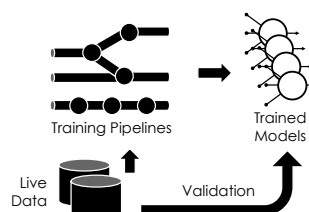
25

## Model Development

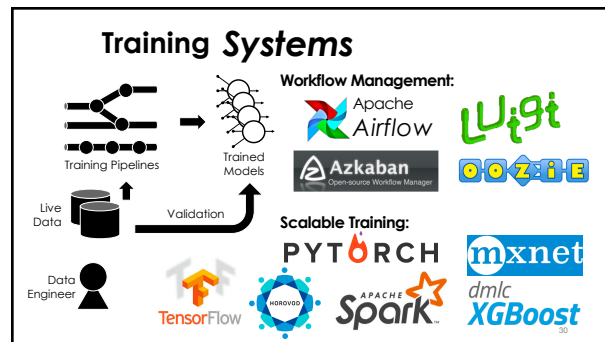


28

## Training

Training models **at scale** on **live data****Retraining** on new dataAutomatically **validate** prediction accuracyManage model **versioning**Requires **minimal expertise** in machine learning

29



30

### Model Training → Hardware

- Fewer models to train → need **distributed training** of individual models
  - Often train with **more data**
- **Larger models** and **mini-batch sizes**
  - Need larger on-device memory
  - Counter trends → reversible networks, **optimal checkpointing**, ...
- Models and hyperparameters are vetted → focus on **system optimizations**
  - Can **tolerate** some **system error** (quantization and async.)
  - Need adequate stability to meet deadlines
- **Data preparation** is still potentially an issue (as with model dev.)
- Need to deal with **composition** of multiple models

31

### Do Bigger Models Train Faster?

(Preliminary unpublished work.)

- Studying **pre-training** of large **Transformer Models** for NLP task (e.g., BERT)

32

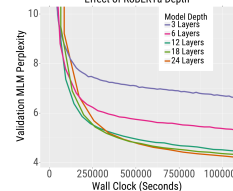
### Do Bigger Models Train Faster?

(Preliminary unpublished work.)

- Studying **pre-training** of large **Transformer Models** for NLP task (e.g., BERT)

#### Deeper Models **Reduce Error Faster**

Effect of RoBERTa Depth



33

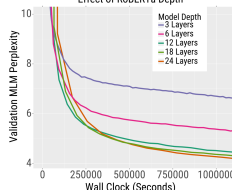
### Do Bigger Models Train Faster?

(Preliminary unpublished work.)

- Studying **pre-training** of large **Transformer Models** for NLP task (e.g., BERT)

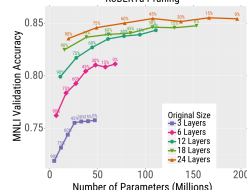
#### Deeper Models **Reduce Error Faster**

Effect of RoBERTa Depth



#### More Resilient to **Lossy Compression**

RoBERTa Pruning

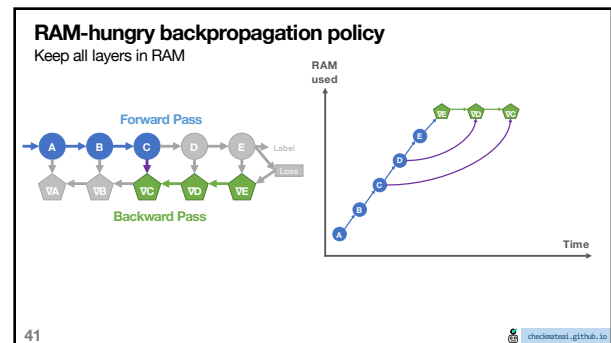
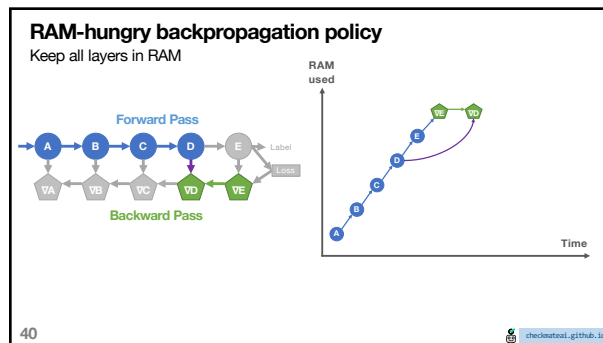
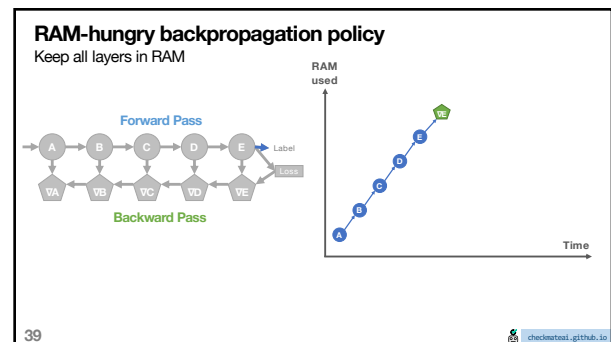
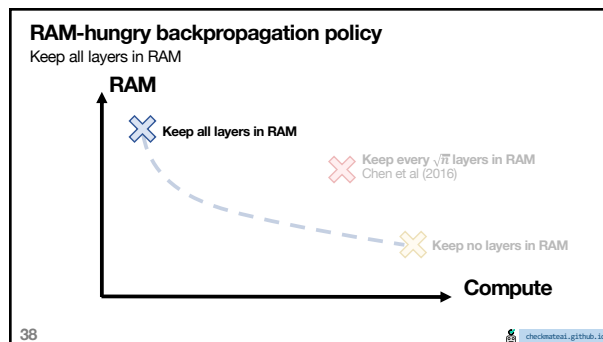
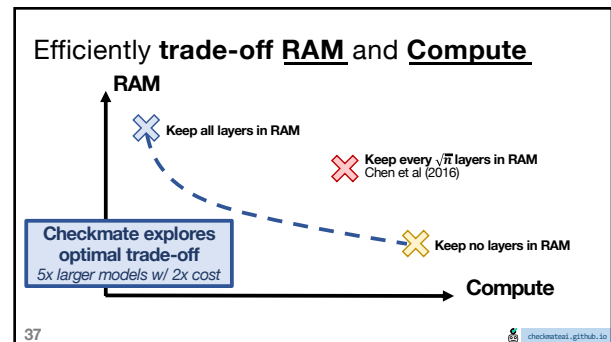
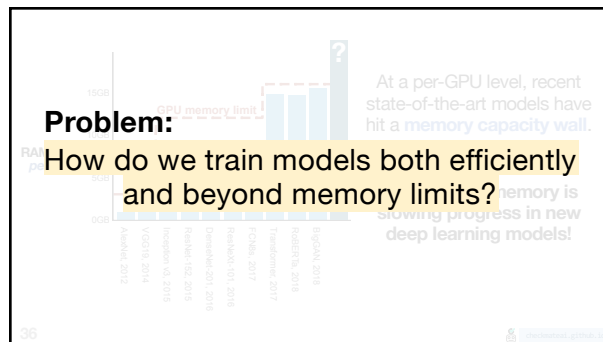


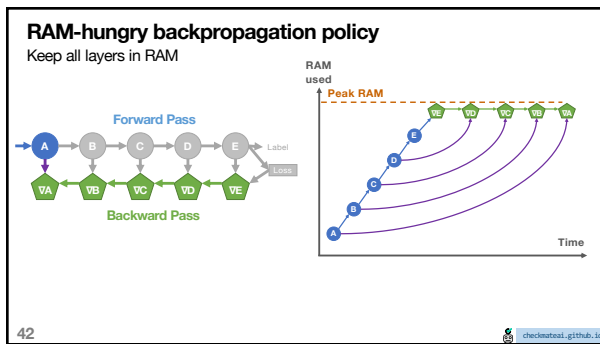
34

### CheckMate [MLSys 2020]

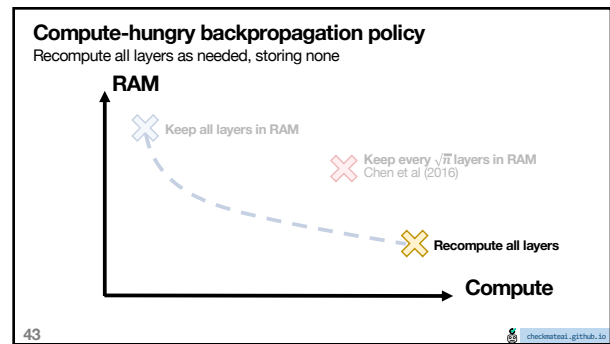
Scaling deep learning training beyond the GPU memory wall.

35

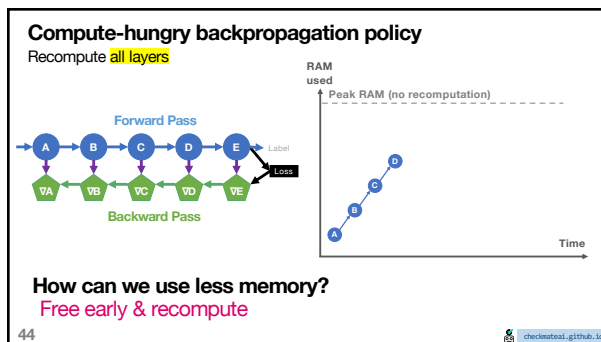




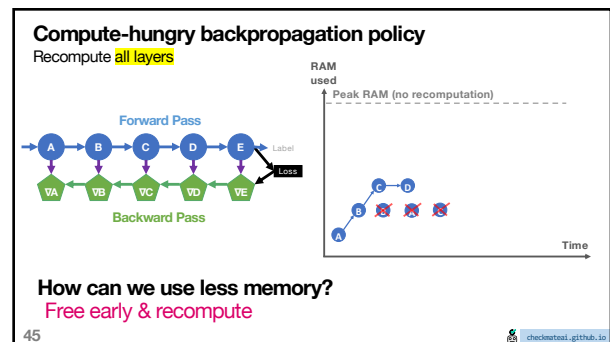
42



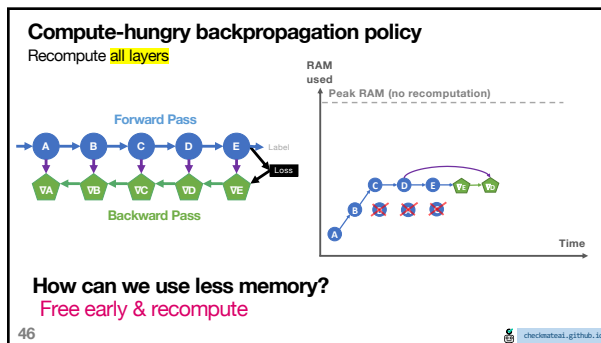
43



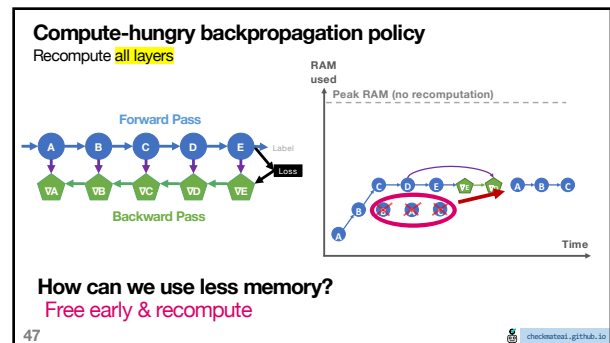
44



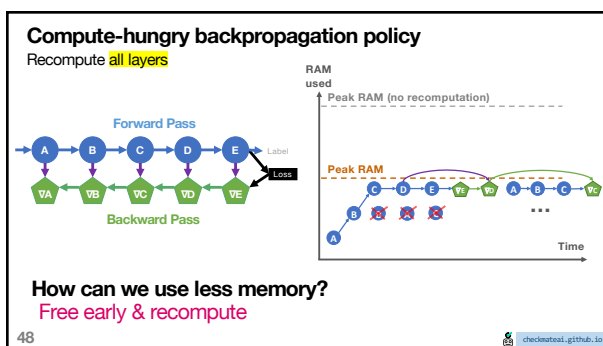
45



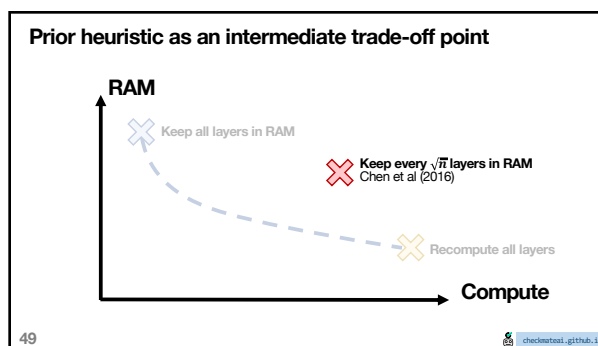
46



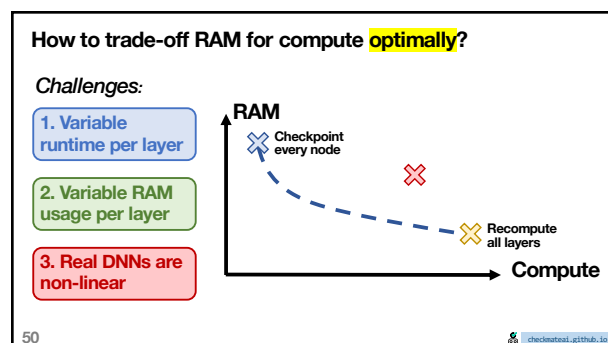
47



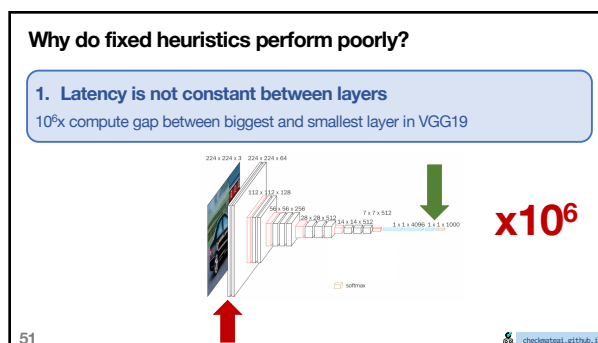
48



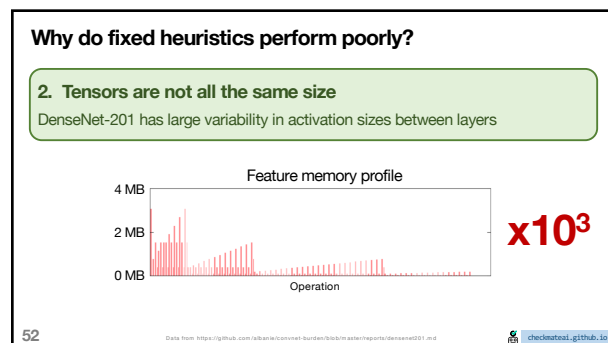
49



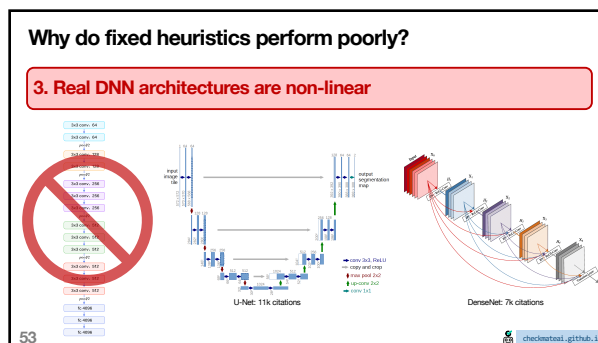
50



51



52



53

## Checkmate

A system for **optimal** tensor rematerialization

- Statically optimize graph once (10s to 1hr)
- Train optimized graph for weeks
- Checkmate composed of 3 parts:
  - Profiling**: hardware/RAM aware schedules
  - Integer LP**: enables finding optimal schedule
  - TF2.0 graph pass**: support TPU, GPU, CPU

**Profile layers**  
↓  
**Solve integer LP**  
↓  
**Rewrite TF2.0 graph**

54 [checkmateai.github.io](https://github.com/checkmateai/checkmate)

54

## Variance Reduction for Quantized Training

(Preliminary unpublished work.)  
Jianfei Chen (Postdoc)

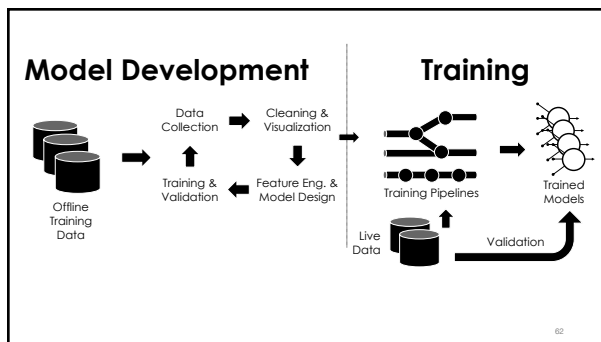
- Quantized weights, activations, and **gradients**
- Studying uniform **stochastic rounding**
  - Prove** gradient is **unbiased**
  - hardware support needed...

**Consequence:**  
Quantization preserves "correctness" of SGD  
→ Does affect convergence rate

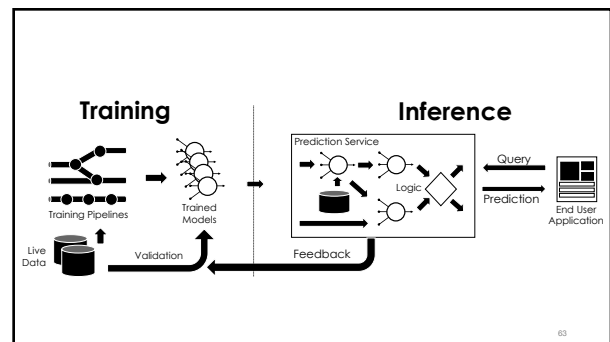
- Analyze Variance**
  - 1 less bit > 2x increase in stdev of gradient estimates
  - 4x increase in batch size to recover convergence rates
- Studying variance reduction mechanism
  - developing **quantization preconditioners**

55

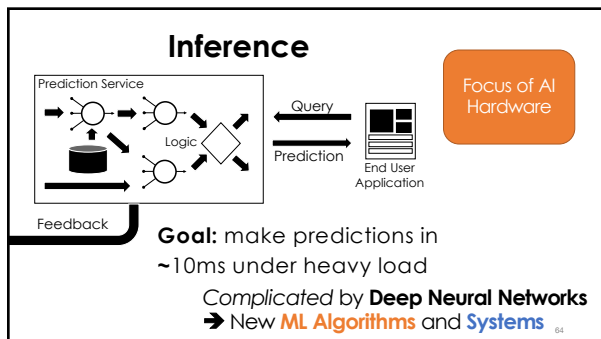
55



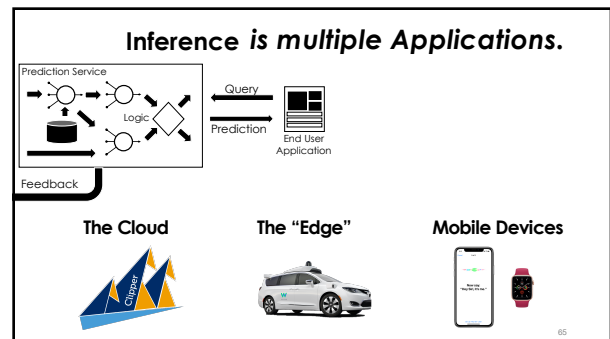
62



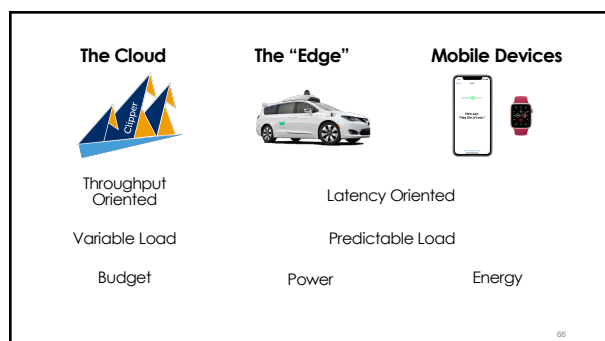
63



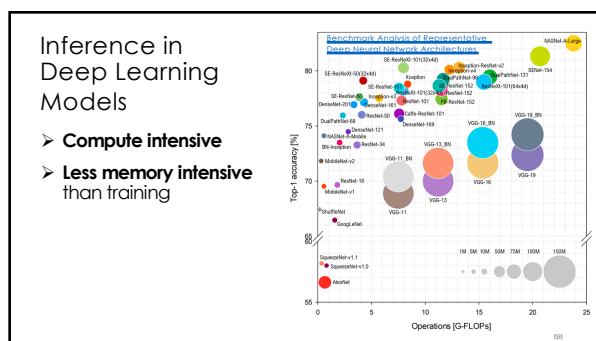
64



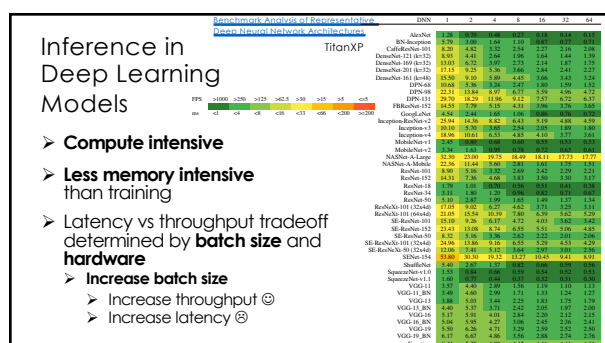
65



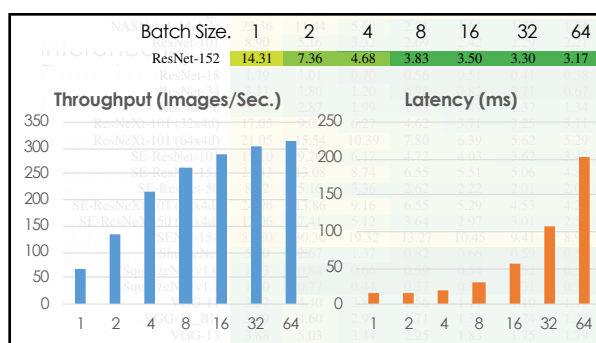
66



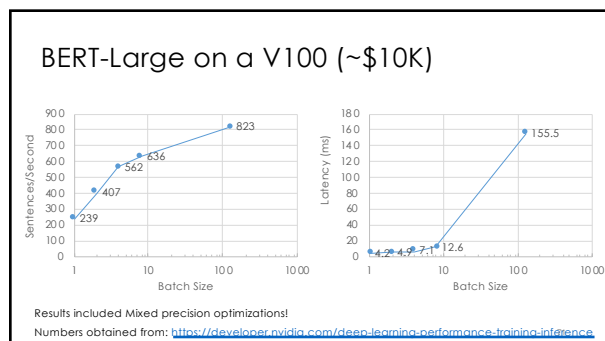
68



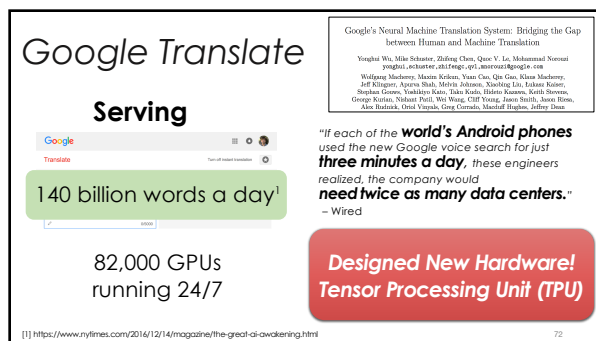
69



70



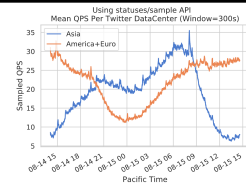
71



72

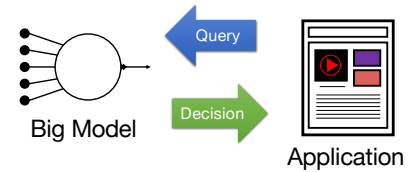
## Other Challenges?

- **Bursty load** →
  - overprovision resources →
  - **expensive**
  - TPU reports 28% utilization of vector units in production
- **Solutions**
  - statistical multiplexing → hardware not designed for multitennancy
  - could try to predict arrival process → generally difficult to predict
- **Versioning and testing models**
- **Prediction pipelines** → more on this soon



73

## Inference

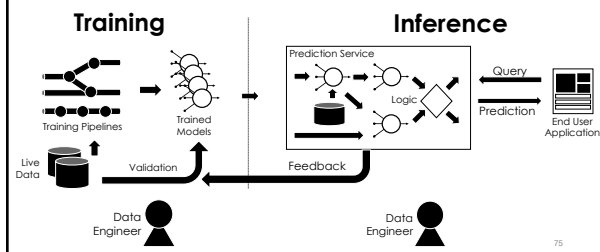


## Two Approaches

- **Offline:** Pre-Materialize Predictions
- **Online:** Compute Predictions on the fly

74

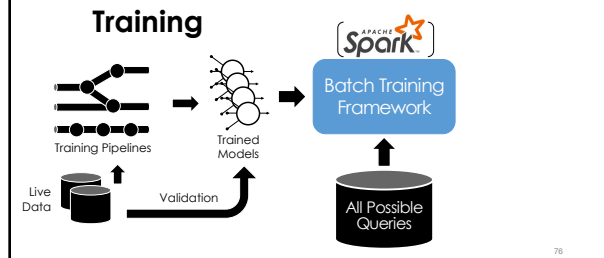
## Pre-materialized Predictions



75

75

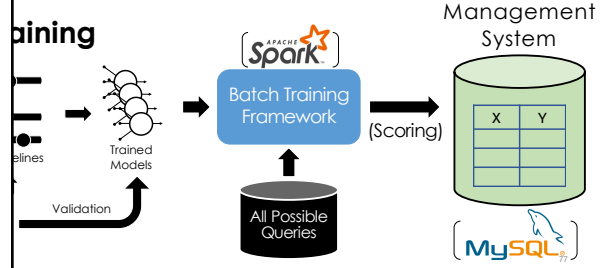
## Pre-materialized Predictions



76

76

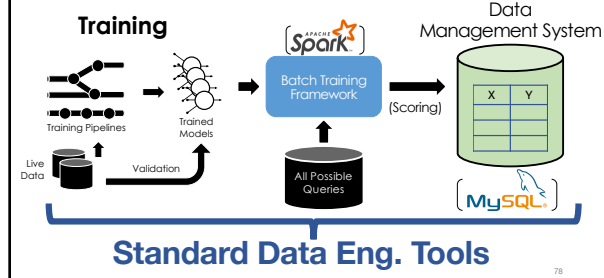
## Pre-materialized Predictions



77

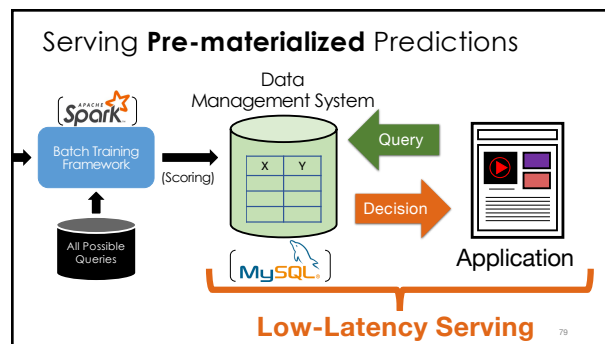
77

## Pre-materialized Predictions

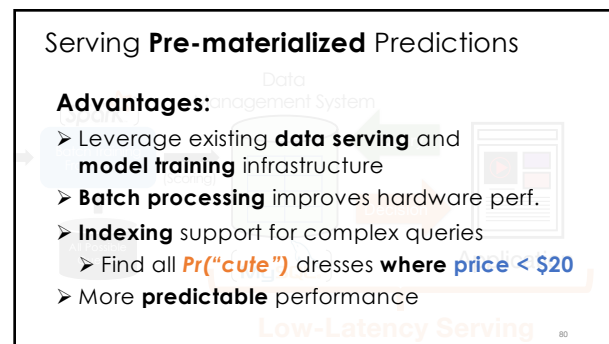


78

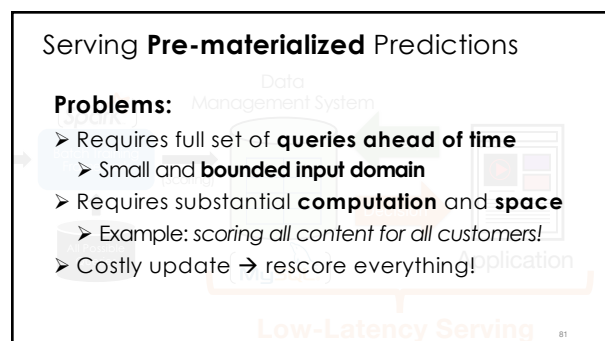
78



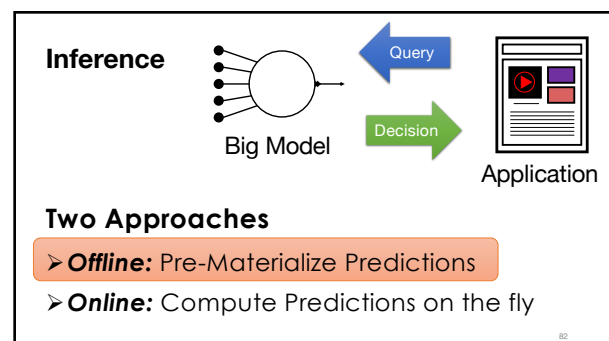
79



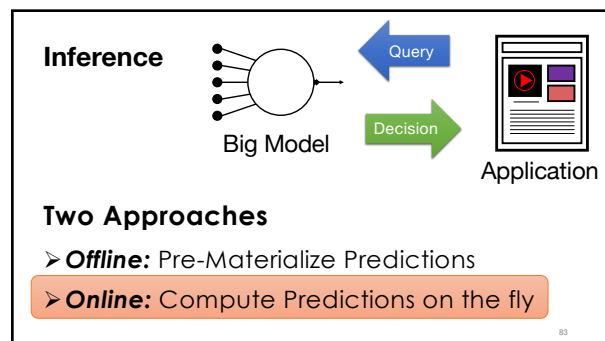
80



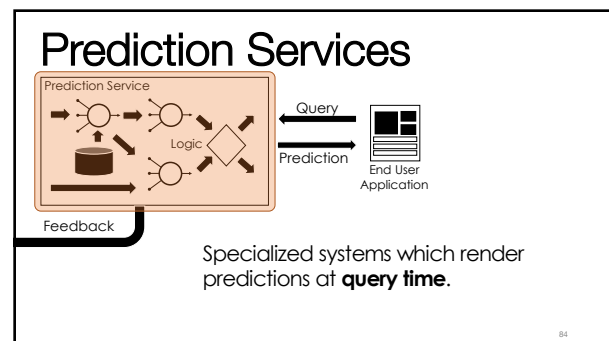
81



82

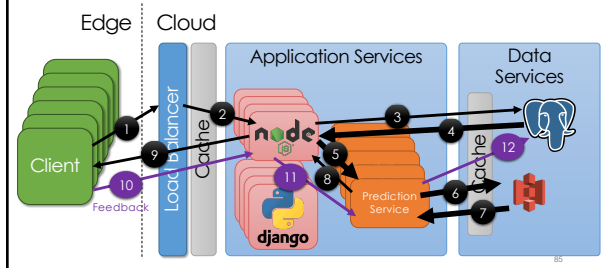


83



84

## Architecture of a Prediction Service



85

## Online: Compute Predictions at Query Time

### Examples

- Signals processing: speech recognition & image tagging
- Ad-targeting based on search terms, available ads, user features

### Advantages

- Compute only necessary queries
- Enables models to be changed rapidly (e.g., bandit exploration)
- Queries do not need to be from small ground set

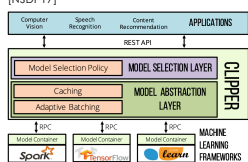
### Disadvantages

- Increases complexity and computation overhead of serving system
- Requires low and predictable latency from models

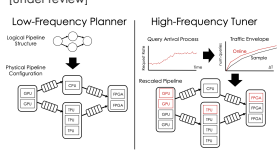
86

## Active Area of Research in my Group

### Clipper Prediction Serving System [NSDI'17]



### InferLine Pipeline Provisioning System [Under review]



87

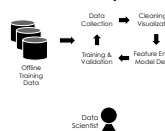
## Prediction Serving → Hardware

- Inference requires less memory → **focus on compute**
- Greater emphasis on **latency** instead of throughput
  - Focus on **small batch** inference (batch size = 1)
  - Opportunity to exploit **pipeline parallelism**
  - Need **high availability** → esp. in mission critical settings
- Often runs multiple **concurrent prediction tasks**
  - Cloud → Multitenancy → Performance isolation
  - Edge → supporting multiple data streams
- Tolerate **model compression and quantization**
  - As low as 4-bit activations and weights
- **Bursty load**
  - Statistical multiplexing
  - Use inference hardware for background training?

88

## Machine Learning Lifecycle

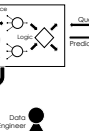
### Model Development



### Training



### Inference

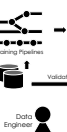


89

### Model Development



### Training



### Inference



# Security?

90

## Protect the **data**, the **model**, and the **query**

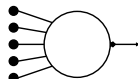
High-Value **Data is Sensitive**

- Medical Info.
- Home video
- Finance



**Models** capture **value** in data

- Core Asset
- “Contain” the data



Biggest opportunity for hardware in ML

**Queries** can be as sensitive as the data

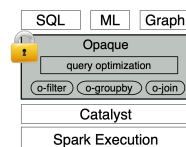


91

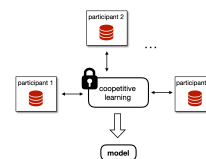
91

## Our recent work in secure ML

**Opaque** | Oblivious Spark over **SGX**  
[NSDI'17]



**Helen** | **Coopetitive Learning**  
Using Cryptographic Primitives  
[S&P'19]



92

92

## Security and Hardware

- Improved Access to Data
  - User willing to share data with models but **not companies (people)**
  - **Differential Privacy** can increase data sharing incentives
- Better **isolation** of co-tenant models on hardware accelerators
- **Coopetitive Learning**: Secure multiparty computation for ML
  - **Example**: Competing banks collaborate to construct a shared fraud model without sharing data.
- Models have access to **more sensitive inputs**
  - **Example**: Alexa could see where you are when asking to turn on a light.

93

93

## Conclusion

- **History**: AI and Computer Systems have Co-evolved
- **Feedback Cycle**: Hardware, Abstractions, and Data
- **ML is many Applications**: *Machine Learning Lifecycle*
  - *Model Development*: Exploration
  - *Training*: Scale and Composition
  - *Inference*: Cloud – Edge Spectrum
- **Security**: Opportunity for hardware innovation in AI

Thank you!

94

94