



Hardware Opportunities in the Machine Learning Lifecycle

Joseph E. Gonzalez
Co-director of the RISE Lab
jegonzal@cs.berkeley.edu

1

1

Get the latest slides and links to literature



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

2

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

3

Outline

- History and the **Co-evolution** of Hardware and AI
- The **Feedback Cycle** driving the 3rd wave of AI
- Machine Learning is **not a single workload**
 - Stages of the **Machine Learning Lifecycle**
- Security** is the key to better accuracy

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

4

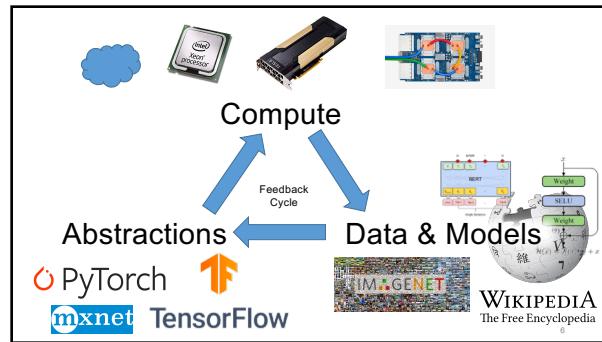
4

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**
 - 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

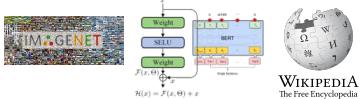
5

5



1

Data & Models



Abstractions



Compute



7

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.

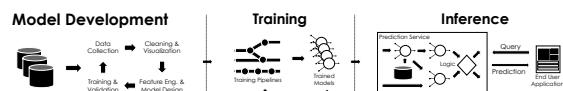
8

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

Machine learning
is not a single application.

9

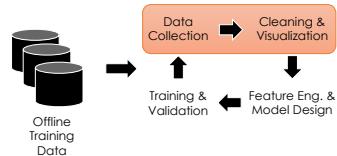
Machine learning is
multiple applications with **different requirements**.



Machine Learning Lifecycle

10

Model Development



Identifying potential sources of data

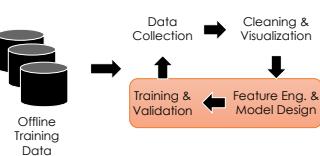
Joining data from multiple sources

Addressing **missing values** and **outliers**

Plotting trends to identify **anomalies**

11

Model Development



Building informative **features functions**

Designing new **model architectures**

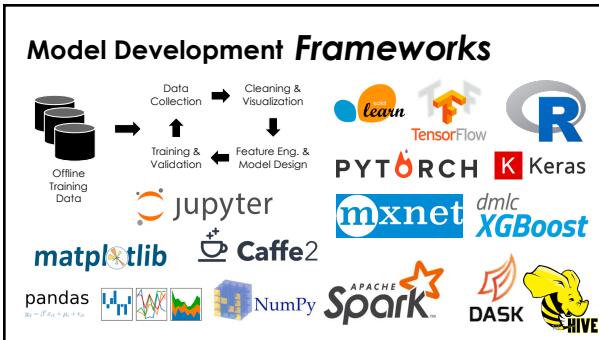
Tuning training algos.

Validating prediction accuracy

12

11

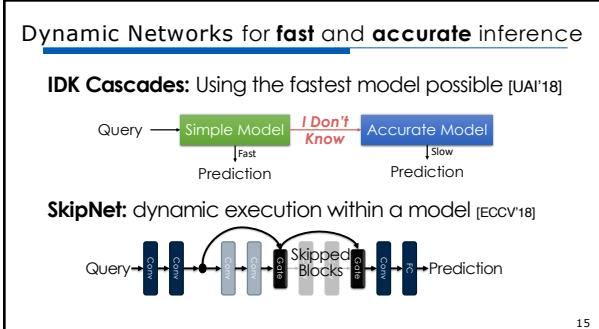
12



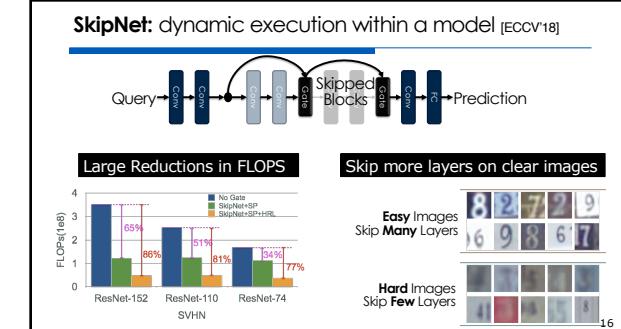
13

- 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** (esp. video)
 - 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

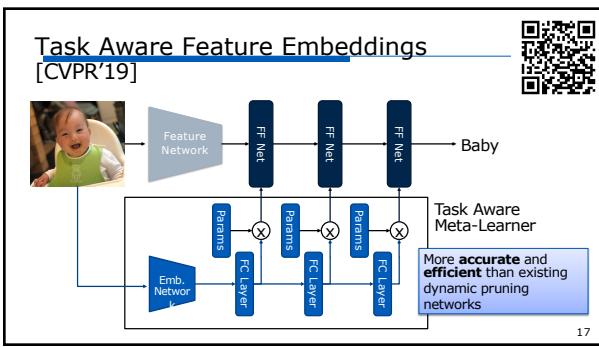
14



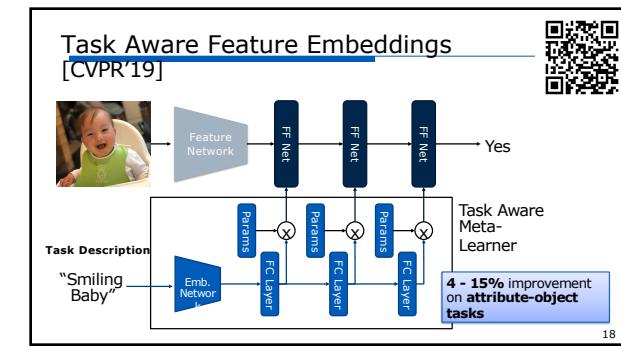
15



16



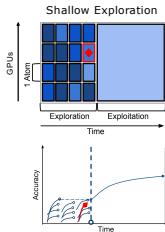
17



18

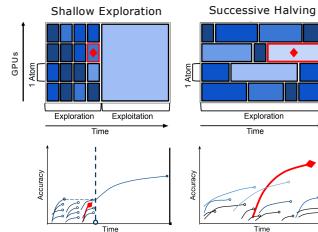
HyperSched [SOCC'19]

Dynamically allocating parallel resources to parallel experiments.



HyperSched [SOCC'19]

Dynamically allocating parallel resources to parallel experiments.

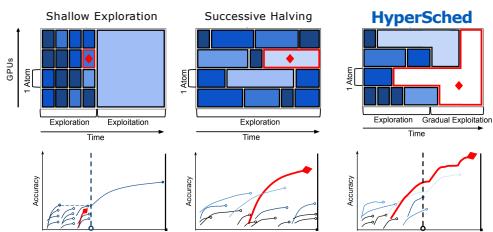


19

20

HyperSched [SOCC'19]

Dynamically allocating parallel resources to parallel experiments.

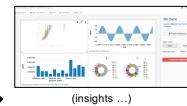


21

What is the output of Model Development



Reports & Dashboards



Trained Model



22

22

What is the output of Model Development



Reports & Dashboards

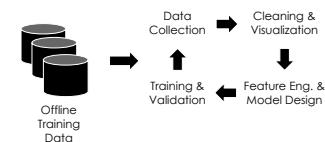


Training Pipelines

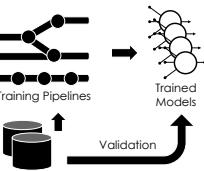


25

Model Development

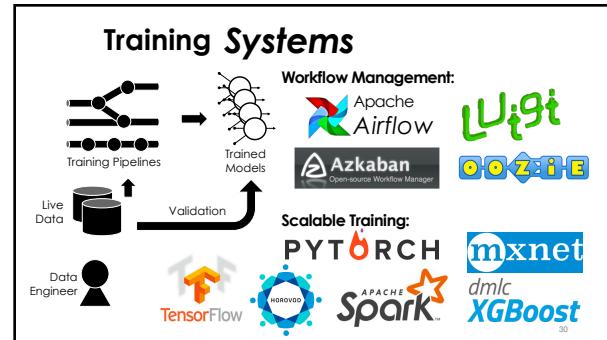
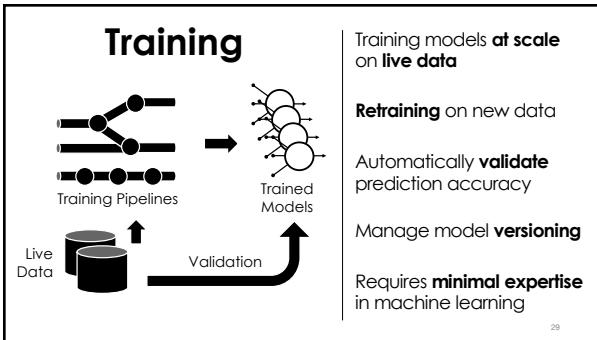


Training



28

28

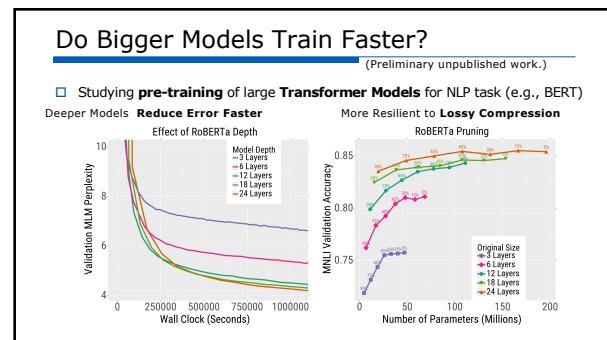


29

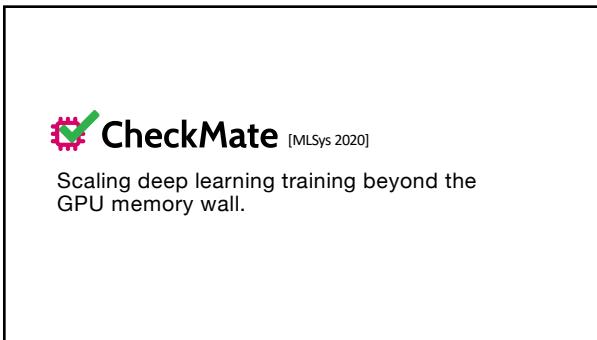
30

- ### Model Training → Hardware
- Fewer models to train → need **distributed training** of **individual models**
 - Often training 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

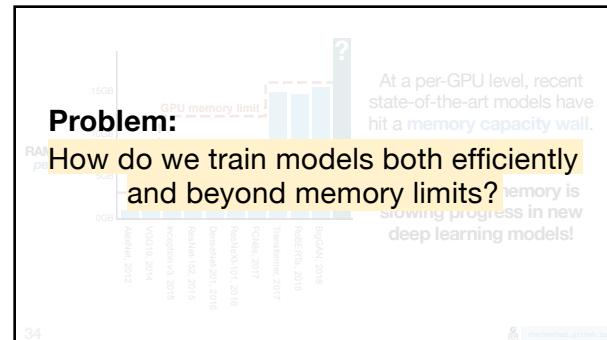
31



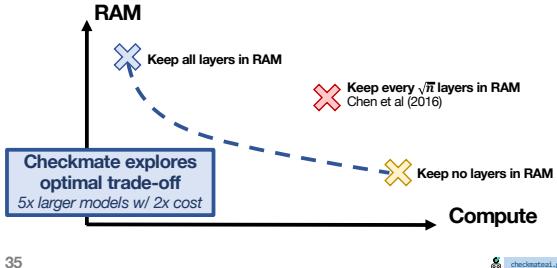
32



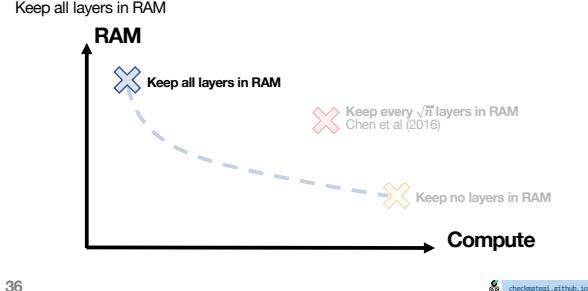
33



Efficiently trade-off RAM and Compute



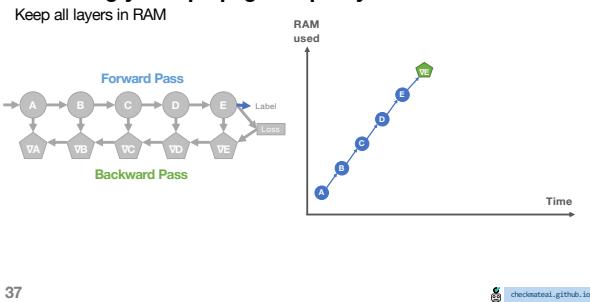
RAM-hungry backpropagation policy



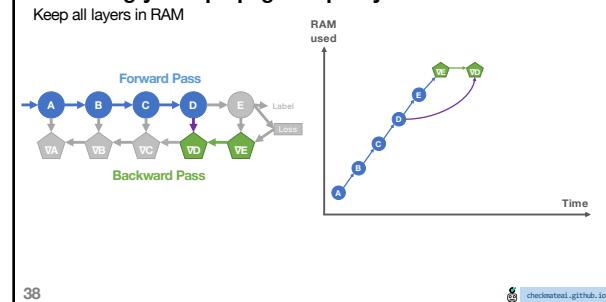
35

36

RAM-hungry backpropagation policy



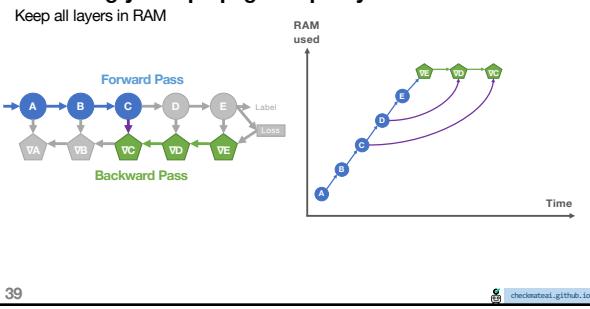
RAM-hungry backpropagation policy



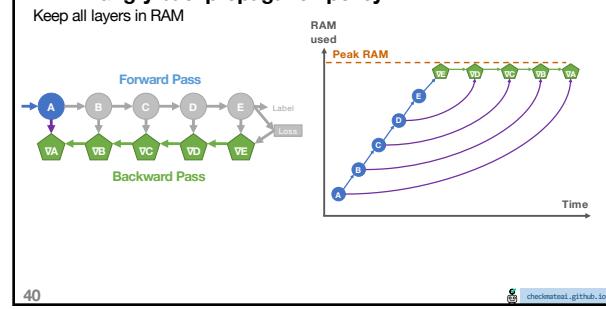
37

38

RAM-hungry backpropagation policy

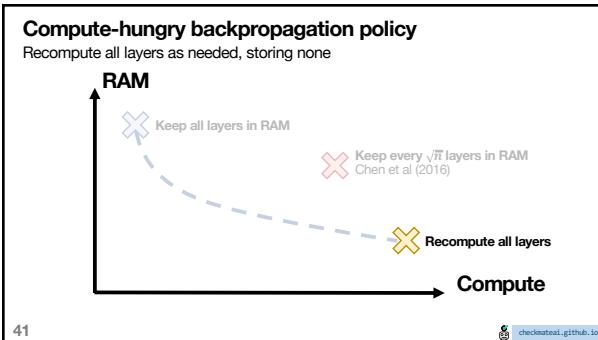


RAM-hungry backpropagation policy

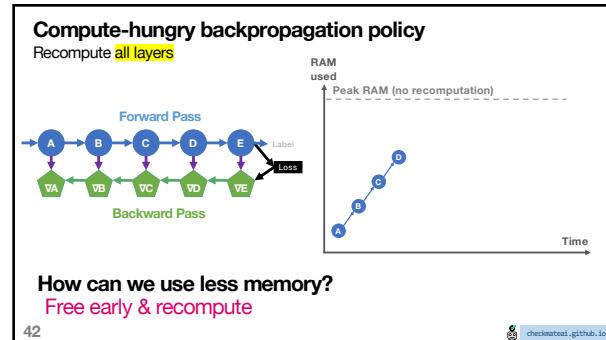


39

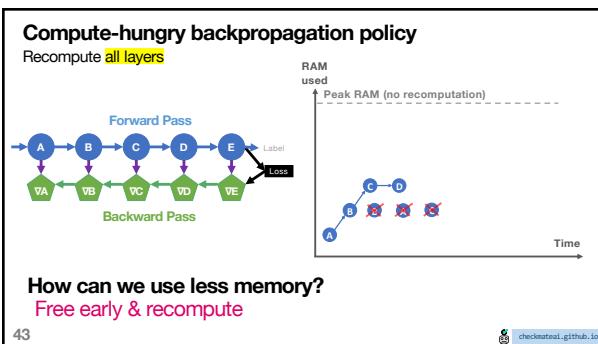
40



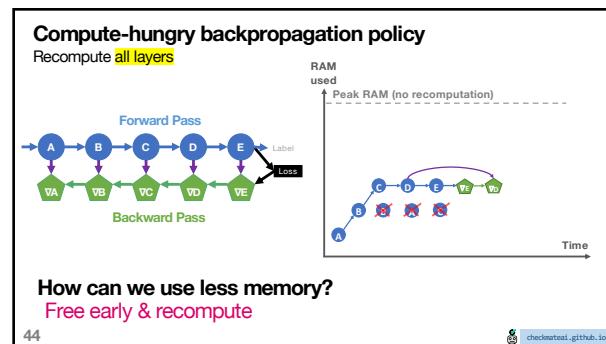
41



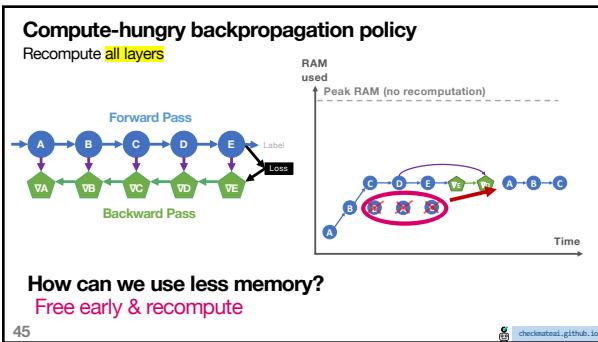
42



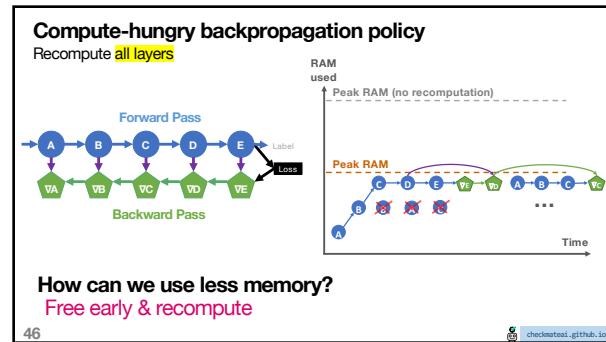
43



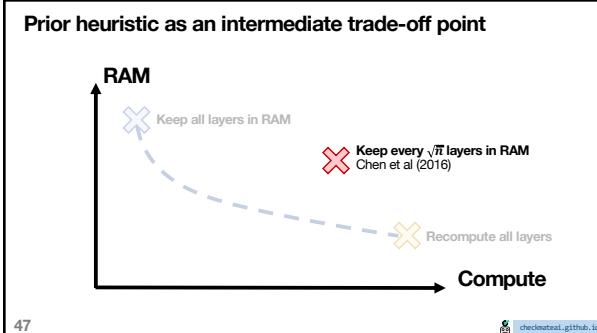
44



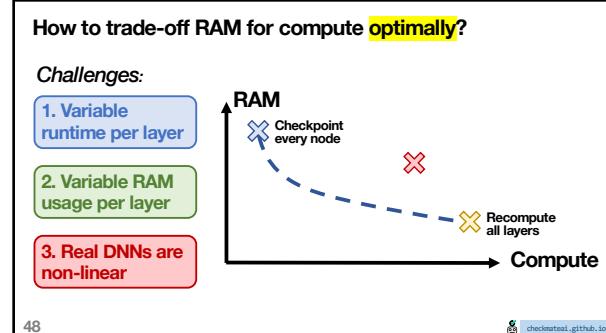
45



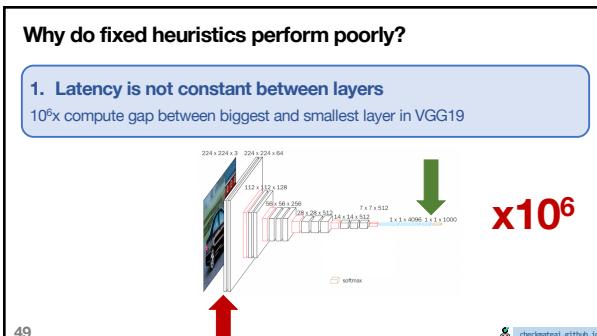
46



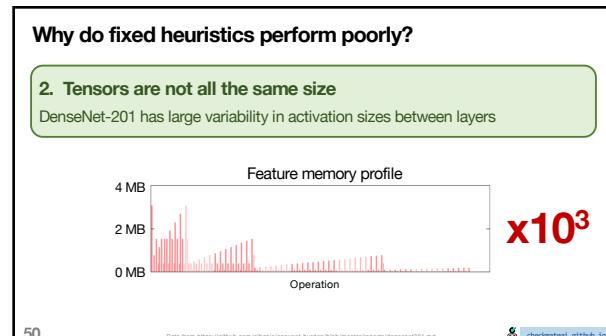
47



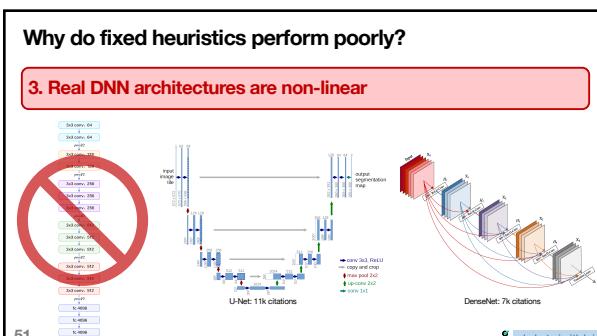
48



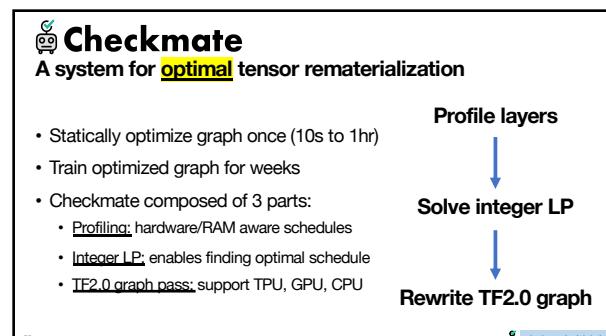
49



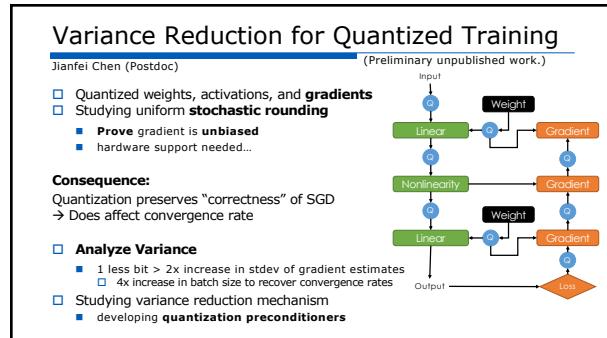
50



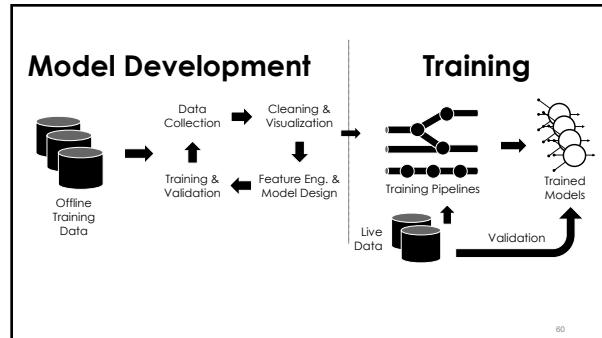
51



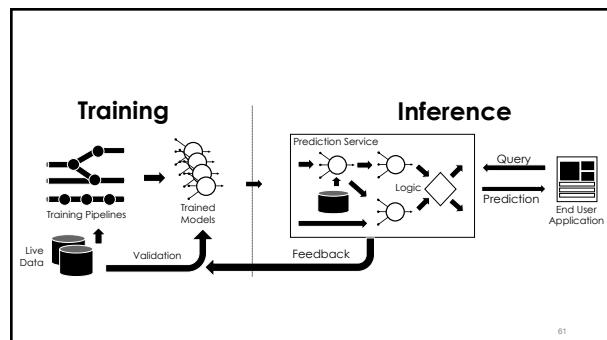
52



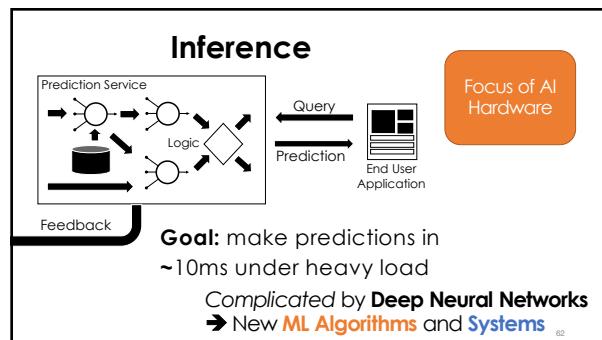
53



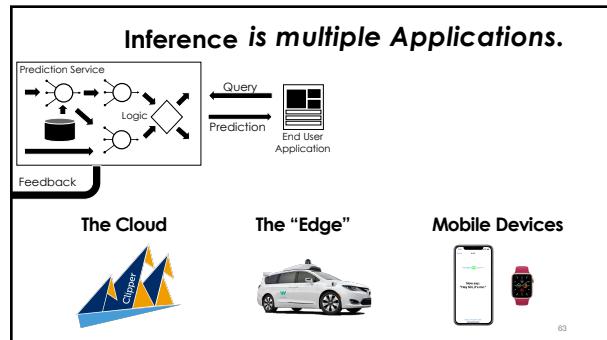
60



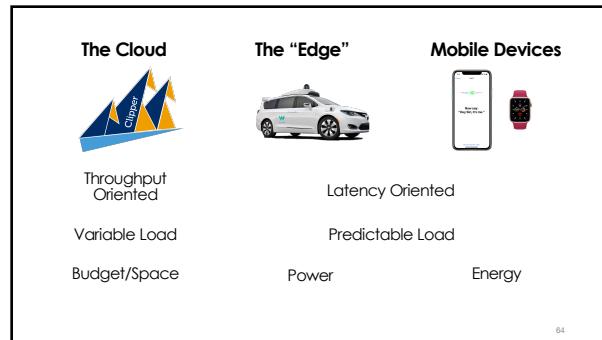
61



62



63



64

Basic Linear Models (Often High Dimensional)

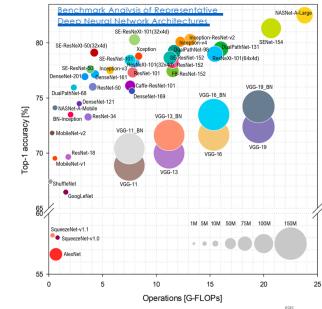
- Common for **click prediction** and **text filter models** (spam)
- Query x encoded in sparse Bag-of-Words:
 - $x = \text{"The quick brown"} = \{(\text{"brown"}, 1), (\text{"the"}, 1), (\text{"quick"}, 1)\}$
- Rendering a prediction: $\text{Predict}(x) = \sigma \left(\sum_{(w,c) \in x} \theta_w c \right)$
- θ is a large vector of weights for each possible word
 - or word combination (n-gram models) ...
- Focus on fast weight retrieval

65

65

Inference in Deep Learning Models

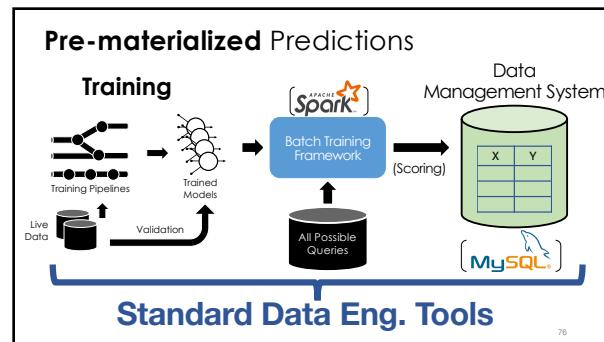
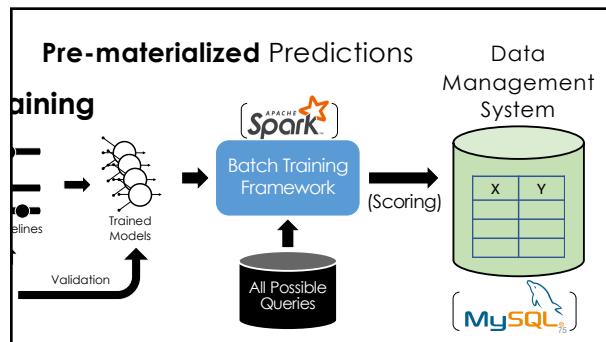
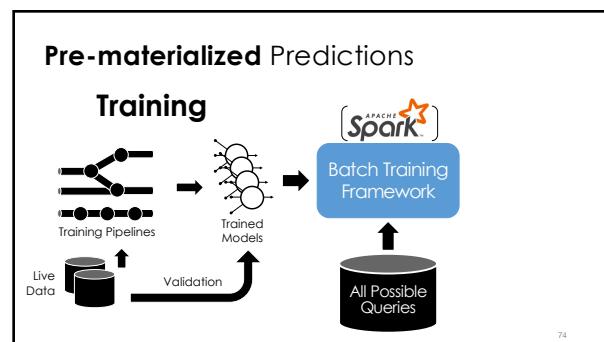
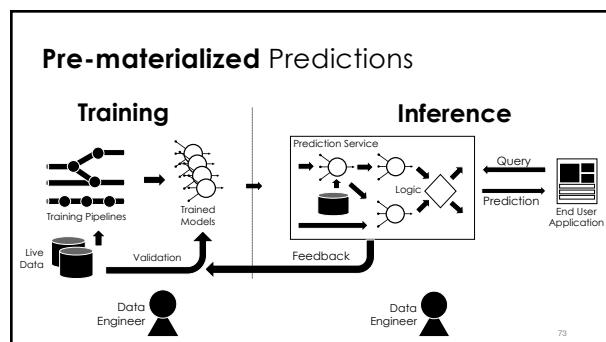
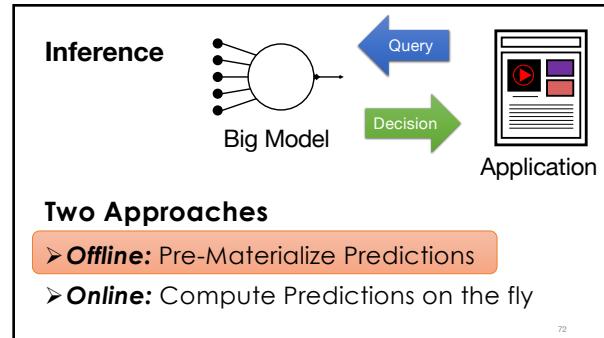
- Compute intensive**
- Less memory intensive** than training

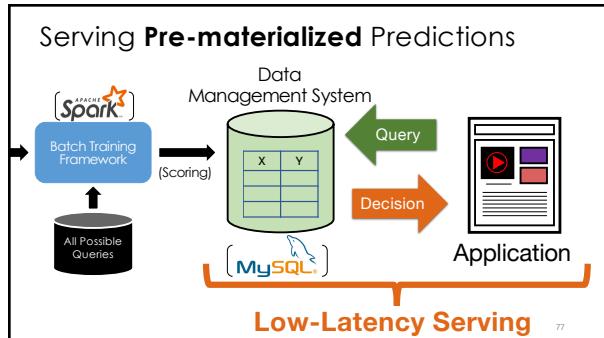


Other Challenges?

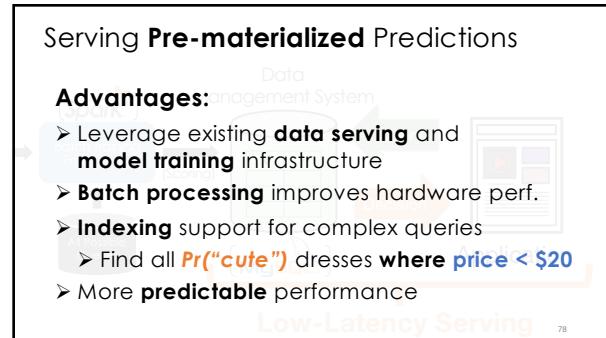
- **Bursty load** →
 - overprovision resources →
 - expensive
 - TPU reports 28% utilization of vector units in production
 - Solutions
 - statistical multiplexing → hardware not designed for multitenancy
 - could try to predict arrival process → generally difficult to predict
- **Versioning and testing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

71

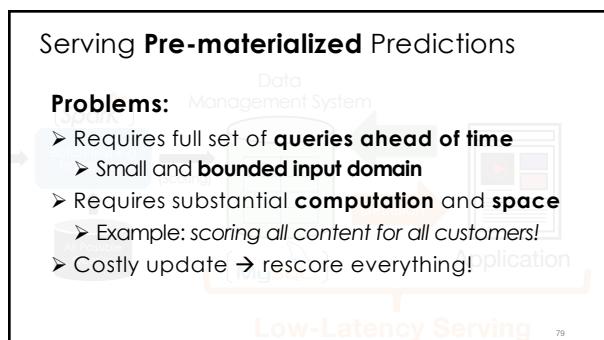




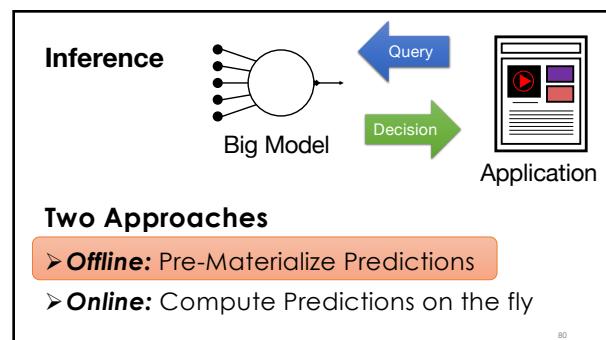
77



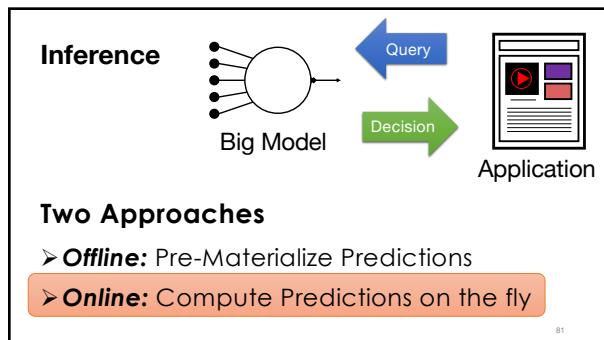
78



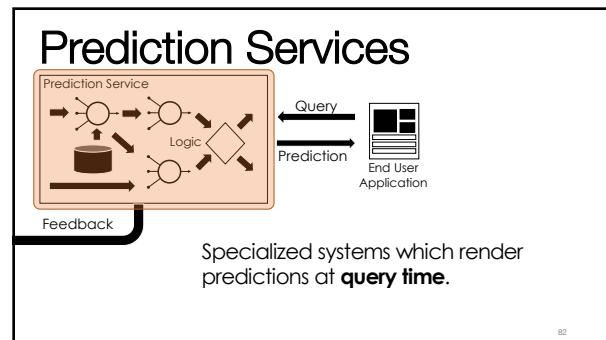
79



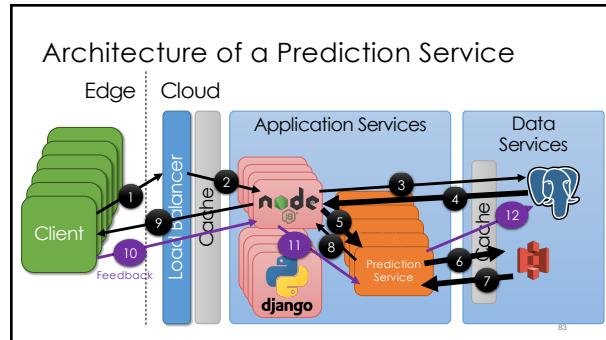
80



81



82

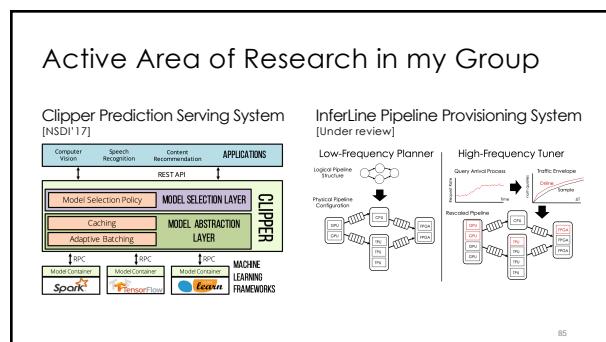


83

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

84

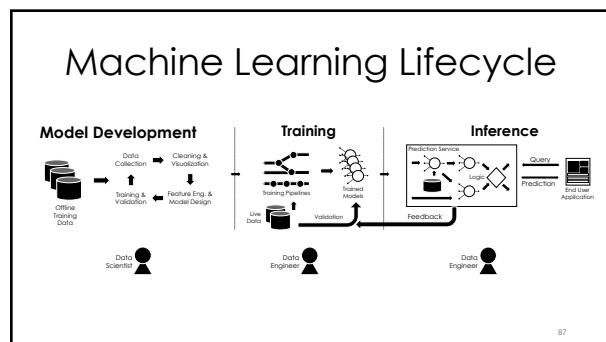


85

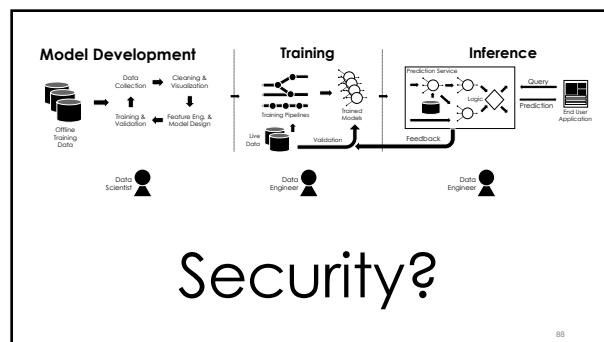
Prediction Serving → Hardware

- 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?

86



87



88

Protect the data, the model, and the query

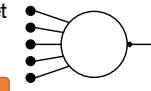
High-Value Data is Sensitive

- Medical Info.
- Home video
- Finance

Biggest opportunity for hardware in ML.



- ### Models capture value in data
- Core Asset
 - “Contain” the data



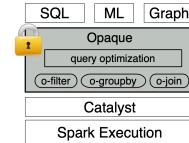
Queries can be as sensitive as the data



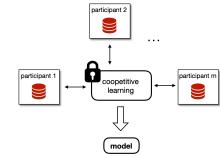
89

Our recent work in secure ML

Opaque | Oblivious Spark over SGX [NSDI'17]



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



90

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
- Cooperative 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.

91

91

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!

92

92