

UC San Diego

# **DSC 102**

# **Systems for Scalable Analytics**

Rod Albuyeh

Topic 5: Model Building Systems

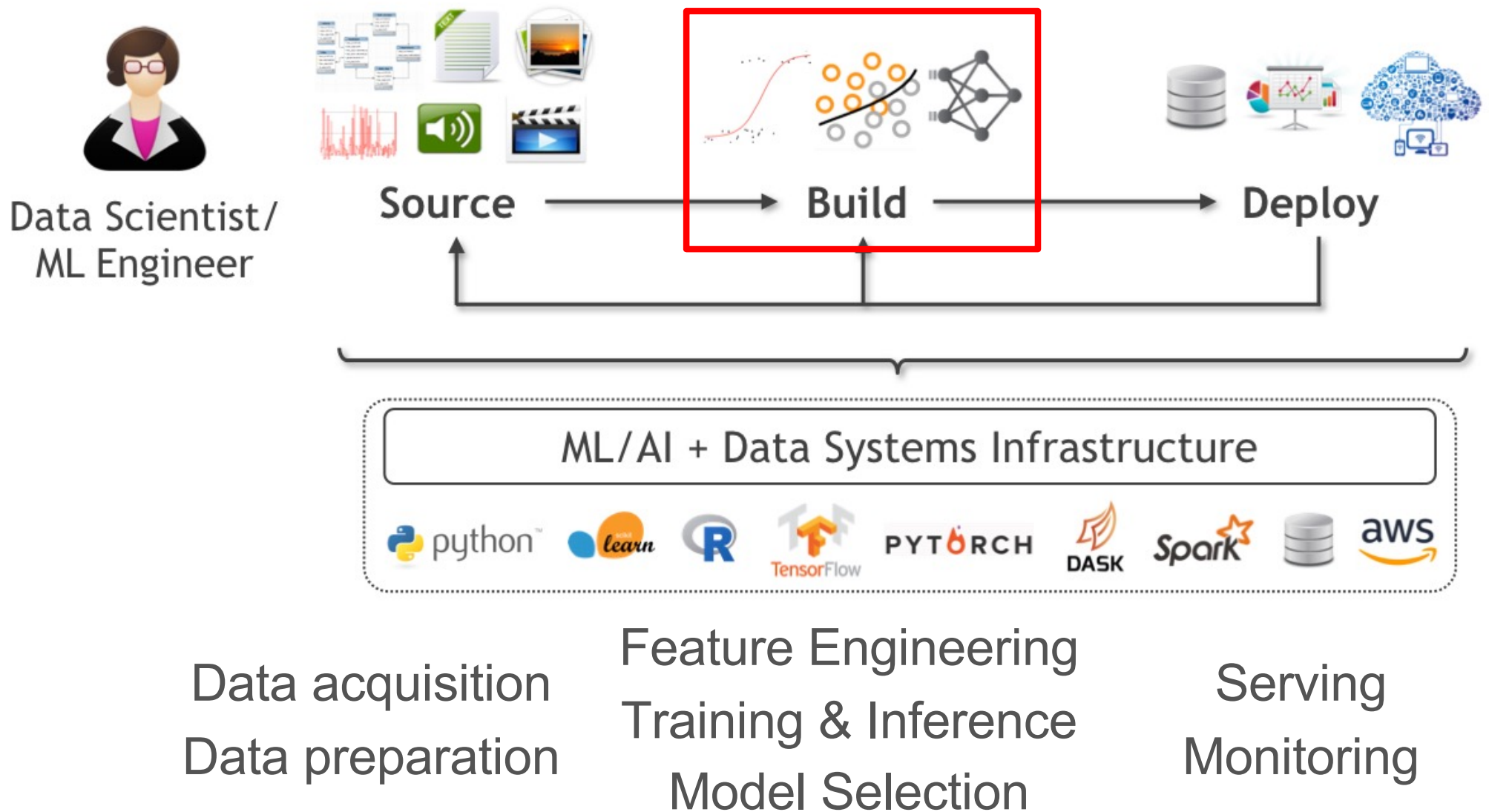
# Admin

Reminder: 90%+ CAPE response rate for class yields 0.5% collective boost to final score.

Current response rate as of June 7<sup>th</sup> 9am is 44.64%.

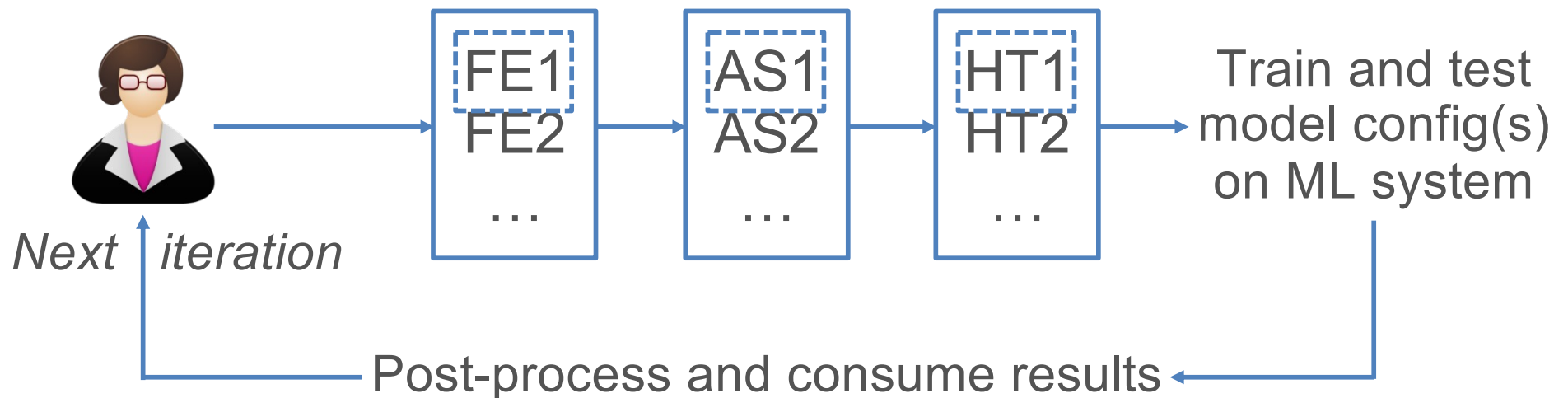
Help both your instructor and DSC 102 improve!

# The Lifecycle of ML-based Analytics

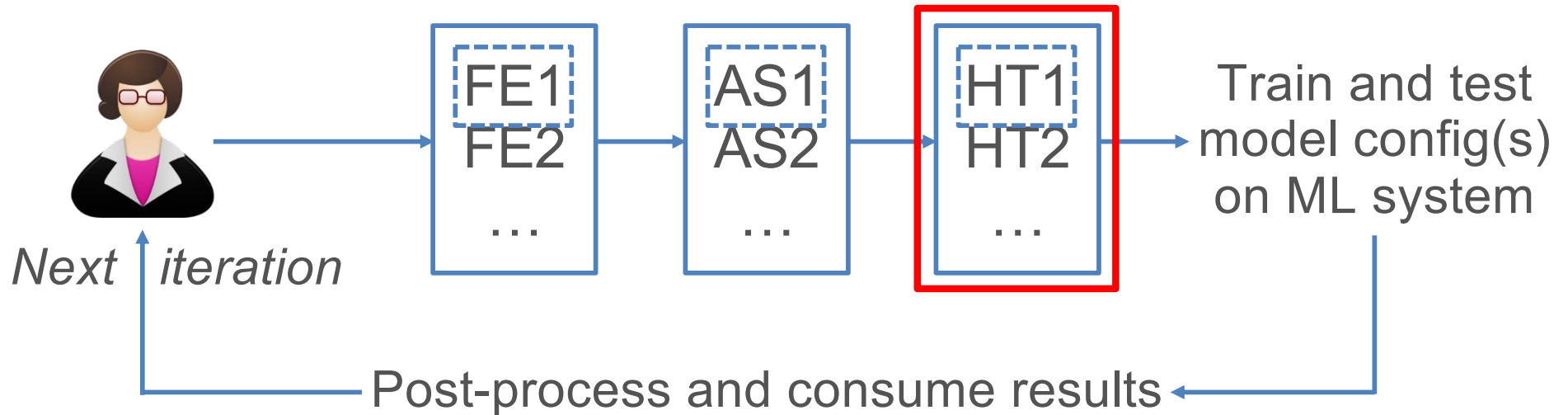


# Model Selection Process

- ❖ Model selection is usually an *iterative exploratory* process with human making decisions on FE, AS, and/or HT
- ❖ Increasingly, automation of some or all parts possible: **AutoML**



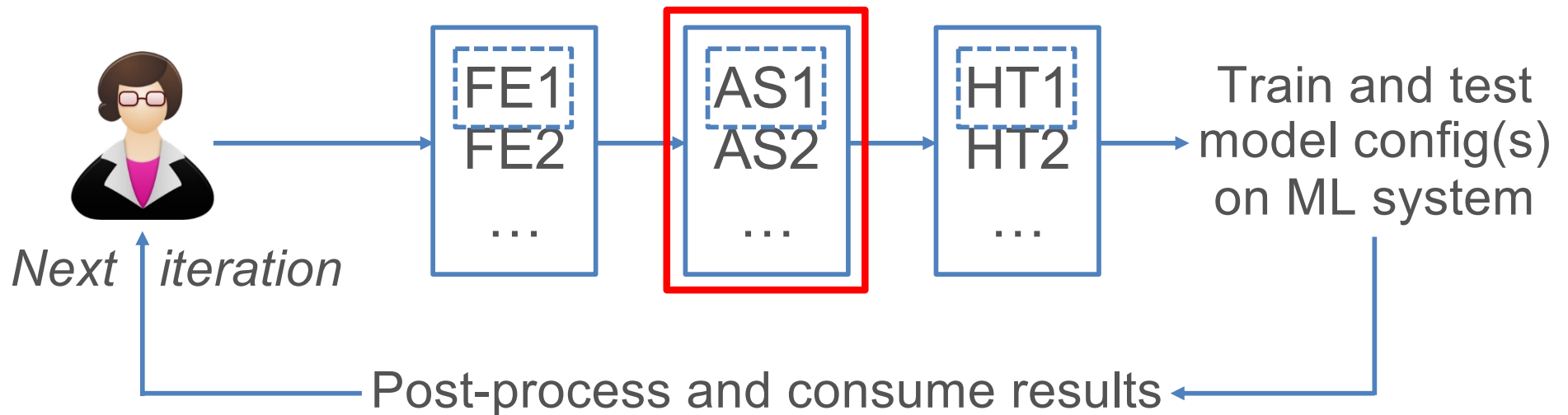
# Hyper-Parameter Tuning



# Hyper-Parameter Tuning

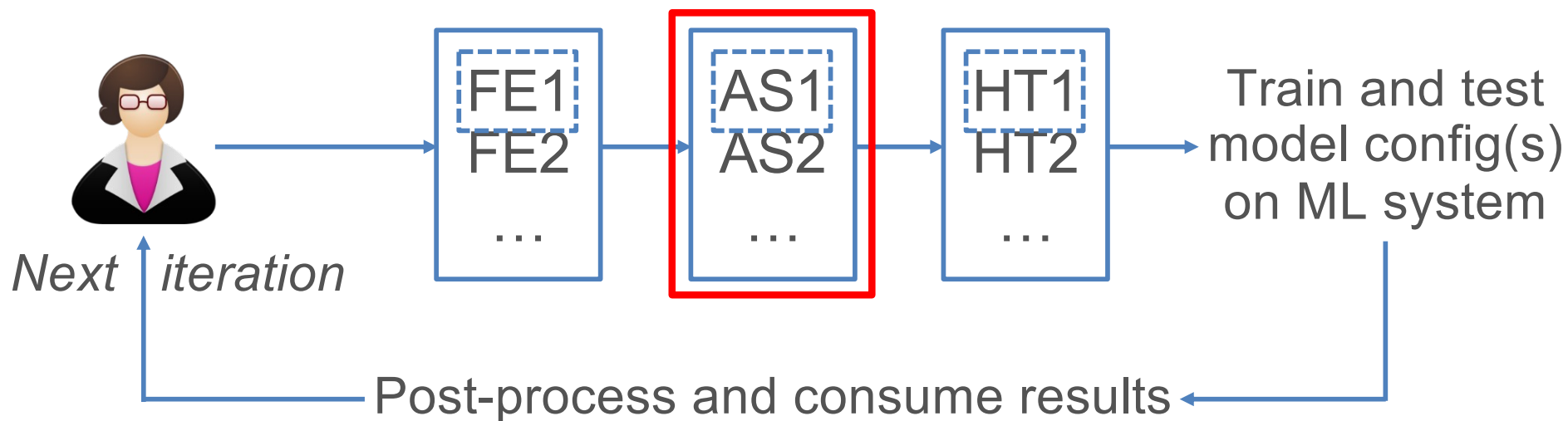
- ❖ **Hyper-parameters:** Knobs for an ML model or training algorithm to control *bias-variance* tradeoff in a dataset-specific manner to make learning effective
- ❖ **Examples:**
  - ❖ GLMs: L1 or L2 *regularizer* to constrain weights
  - ❖ All gradient methods: *learning rate*
  - ❖ Mini-batch Stochastic Gradient Descent: *batch size*
  - ❖ *Others?*
- ❖ HT is an “outer loop” around training/inference
- ❖ Most common approach: **grid search**; pick set of values for each hyper-parameter and take cartesian product
- ❖ Also common: **random search** to subsample from grid
- ❖ Complex AutoML heuristics exist too for HT, e.g., Bayesian

# Algorithm Selection in “classical” ML



- ❖ Not much to say; ML user typically picks models/algorithms in advance
- ❖ Best practice: first train more simple models (log. reg.) as baselines; then try more complex models (XGBoost)
- ❖ **Ensembles:** Build diverse models and aggregate predictions. Even for tabular data, ensembles yield better results and often win Kaggle comps with a few % boost in performance.

# Architecture Selection in DL



- ❖ More critical in DL; neural arch. is **inductive bias** in classical ML parlance; controls feature learning and bias-variance tradeoff
- ❖ Some applications: Many off-the-shelf pre-trained DL models to do “transfer learning,” e.g., see models at [HuggingFace.co](https://huggingface.co)
- ❖ Other applications: Swap pain of hand-crafted feature eng. for pain of neural arch. eng.! Neural arch probably a better interview skill 😊



# Automated Model Selection / AutoML

*Q: Can we automate the whole model selection process?*

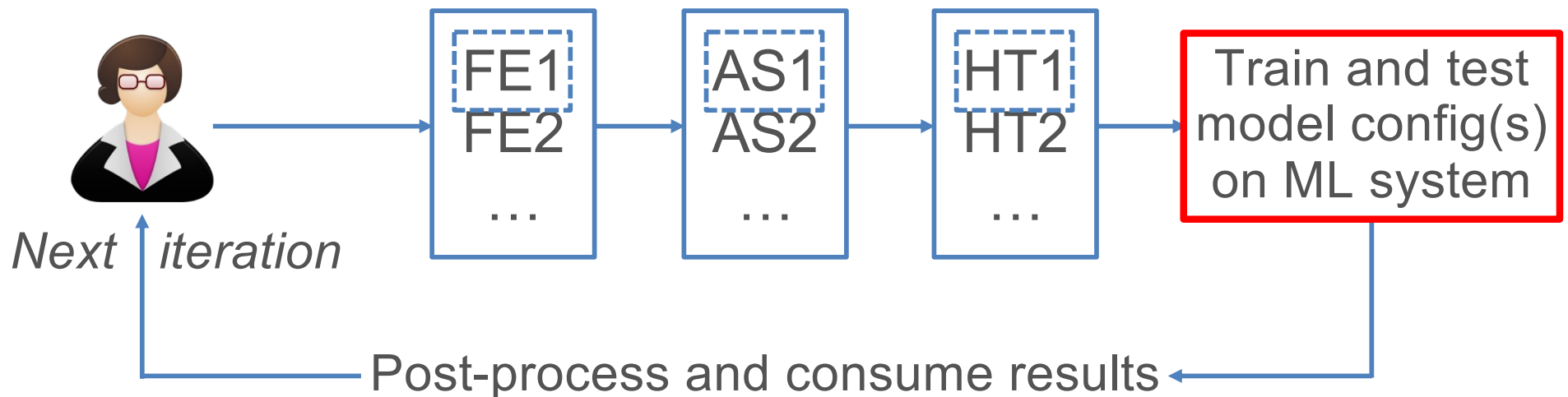
- ❖ It depends. HT and most of FE already automated mostly in practice; (neural) AS is often application-dictated
- ❖ AutoML tools/systems now aim to reduce data scientist's work; or even replace them?! ;)



- ❖ **Pros:** Ease of use; lower human cost; easier to audit; improves ML accessibility
- ❖ **Cons:** Higher resource cost; less user control; may waste domain knowledge; may leave performance on the table
- ❖ Pareto-optima; hybrids possible

**But:** The data sourcing + feature engineering stage is still very hard to automate and tends to be domain / context specific!

# Scalable ML Training and Inference



Then deploy?

# Major ML Model Families/Types

**Generalized Linear Models (GLMs)**; from statistics

**Bayesian Networks**; inspired by causal reasoning

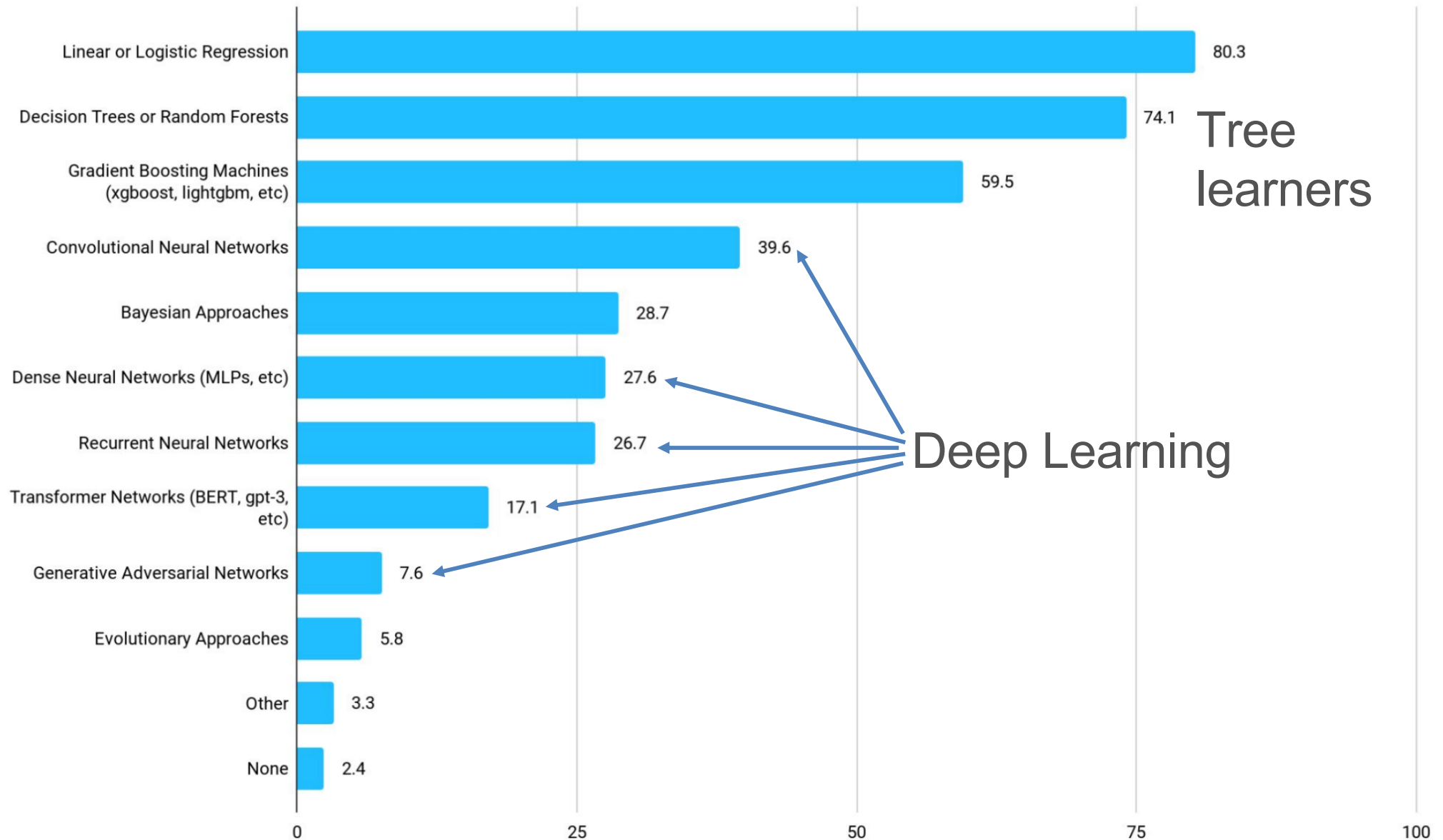
**Decision Tree-based**: CART, Random Forest, Gradient-Boosted Trees (GBT), etc.; inspired by symbolic logic

**Support Vector Machines (SVMs)**; inspired by psychology

**Artificial Neural Networks (ANNs)**: Multi-Layer Perceptrons (MLPs), Convolutional NNs (CNNs), Recurrent NNs (RNNs), Transformers, etc.; inspired by brain neuroscience

**Unsupervised**: Clustering (e.g., K-Means), Matrix Factorization, Latent Dirichlet Allocation (LDA), etc.

# ML Models in Kaggle 2021 Survey



# Scalable ML Training Systems

- ❖ Scaling ML training is involved and model type-dependent
- ❖ Orthogonal Dimensions of Categorization:
  - 1. Scalability:** In-memory libraries vs Scalable ML system (works on larger-than-memory datasets)
  - 2. Target Workloads:** General ML library vs Decision tree-oriented vs Deep learning, etc.
  - 3. Implementation Reuse:** Layered on top of scalable data system vs Custom from-scratch framework

# Major Existing ML Systems

## General ML libraries:

In-memory:



Disk-based files:



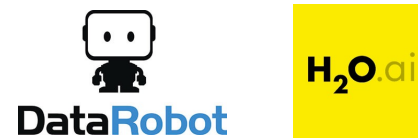
Layered on RDBMS/Spark:



Cloud-native:



“AutoML” platforms:



Decision tree-oriented:



Deep learning-oriented:



# Scalable ML Inference

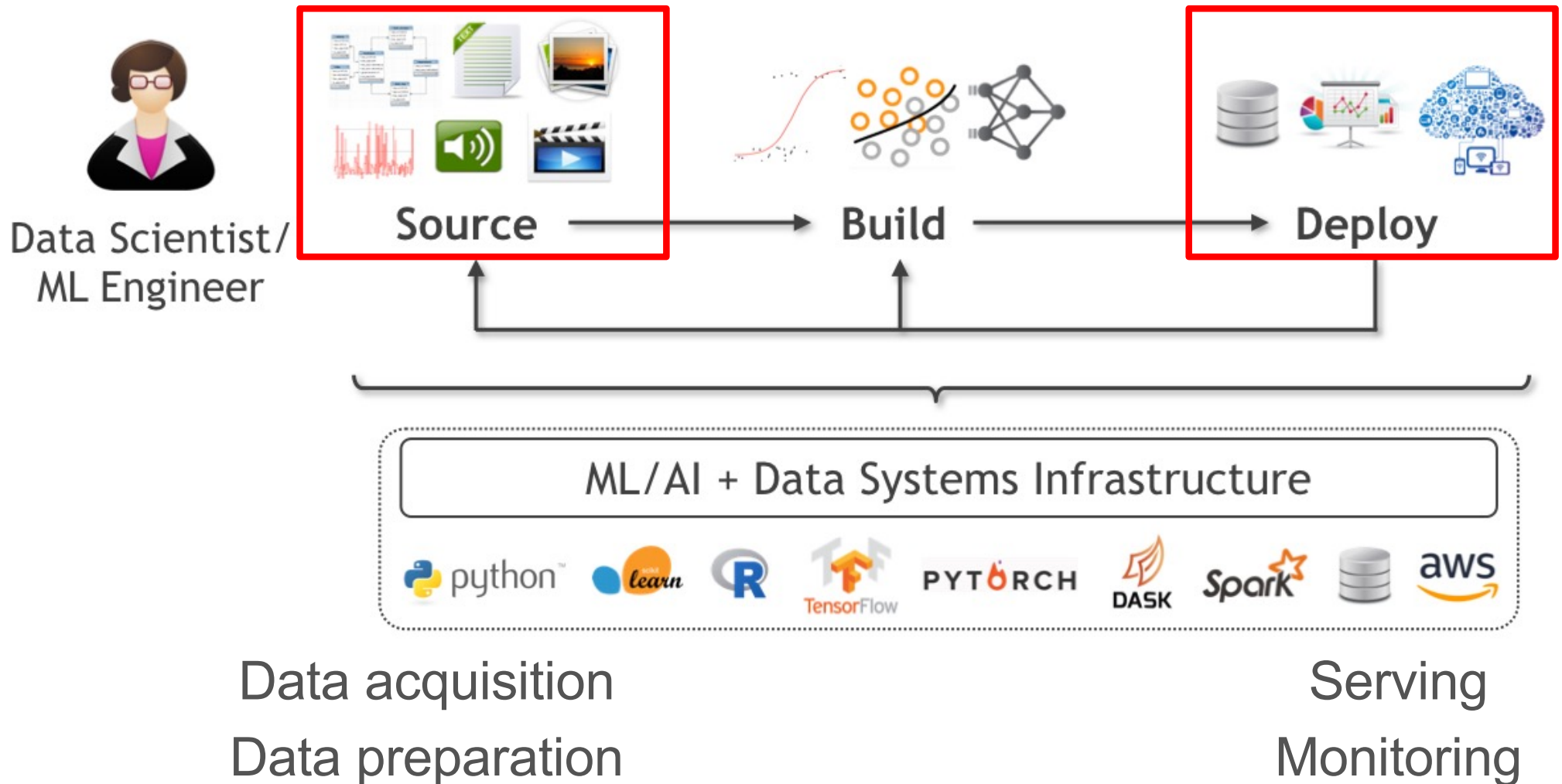
- ❖ A trained/learned ML model is just a prediction function:

$$f : \mathcal{D}_X \rightarrow \mathcal{D}_Y$$

*Q: Given large dataset of examples, how to scale inference?*

- ❖ Assumption 1: An example fits entirely in DRAM
- ❖ Assumption 2:  $f$  fits entirely in DRAM
- ❖ If both hold, trivial access pattern: single filesan, apply per-tuple function  $f$ , write output. How to do this with MapReduce?
- ❖ If either fails, access pattern becomes more complex and dependent on breaking up internals of  $f$  to stage access to data for partial computations

# The Lifecycle of ML-based Analytics



Deployment involves re-applying any transformations and the final  $f()$



# Model Serving / Deployment

- ❖ A trained/learned ML model is just a prediction function:

$$f : \mathcal{D}_X \rightarrow \mathcal{D}_Y$$

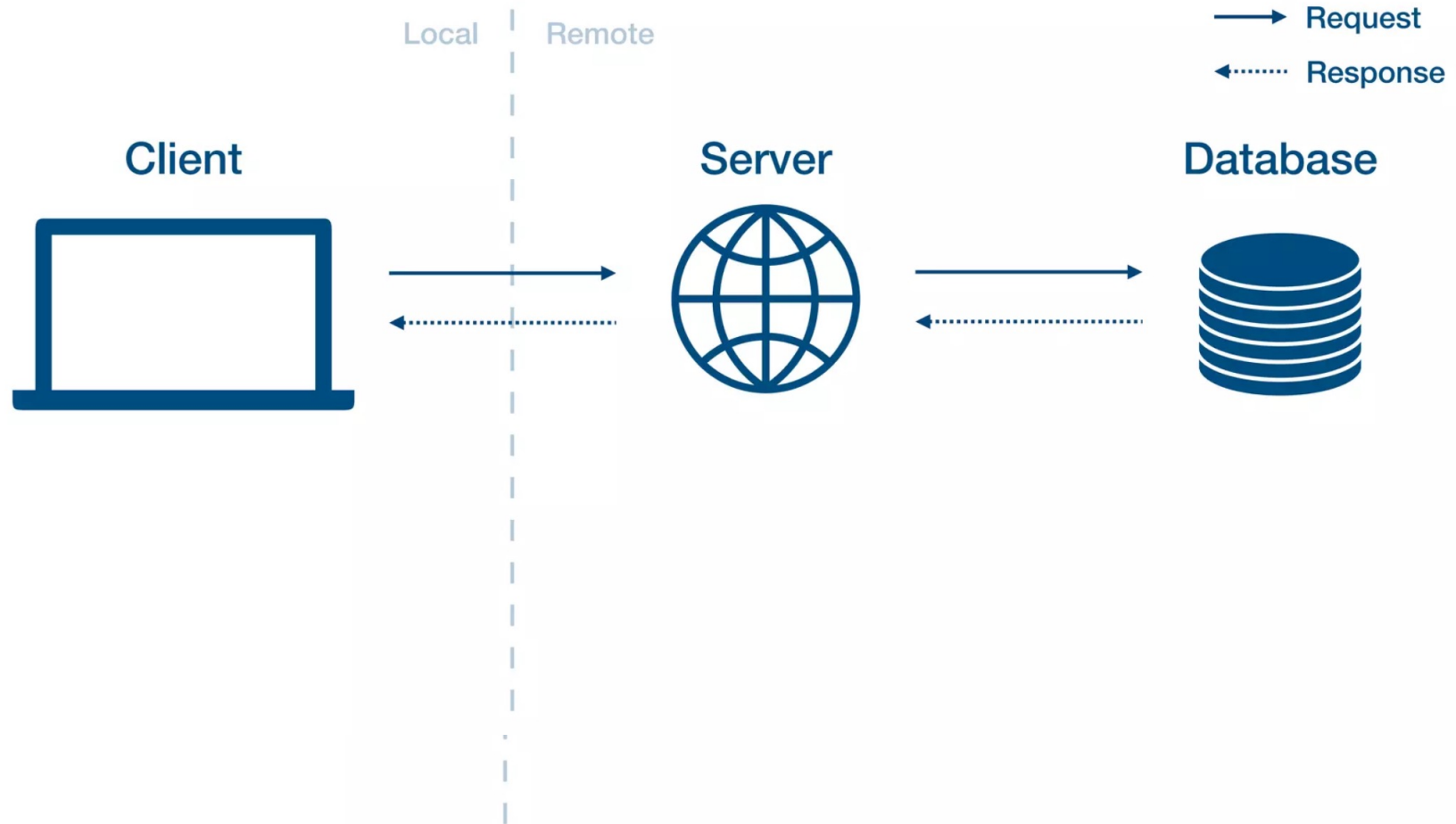
- ❖ A major consideration is, online/realtime vs. offline/batch.
- ❖ In the offline scenario, serving a model is more trivial where it is another processing function that we apply.
- ❖ In the online scenario, we become concerned with millisecond latency for responses, setting up APIs, load balancing, and monitoring.

**DATA SCIENTIST AFTER WRITING A FLASK APP**

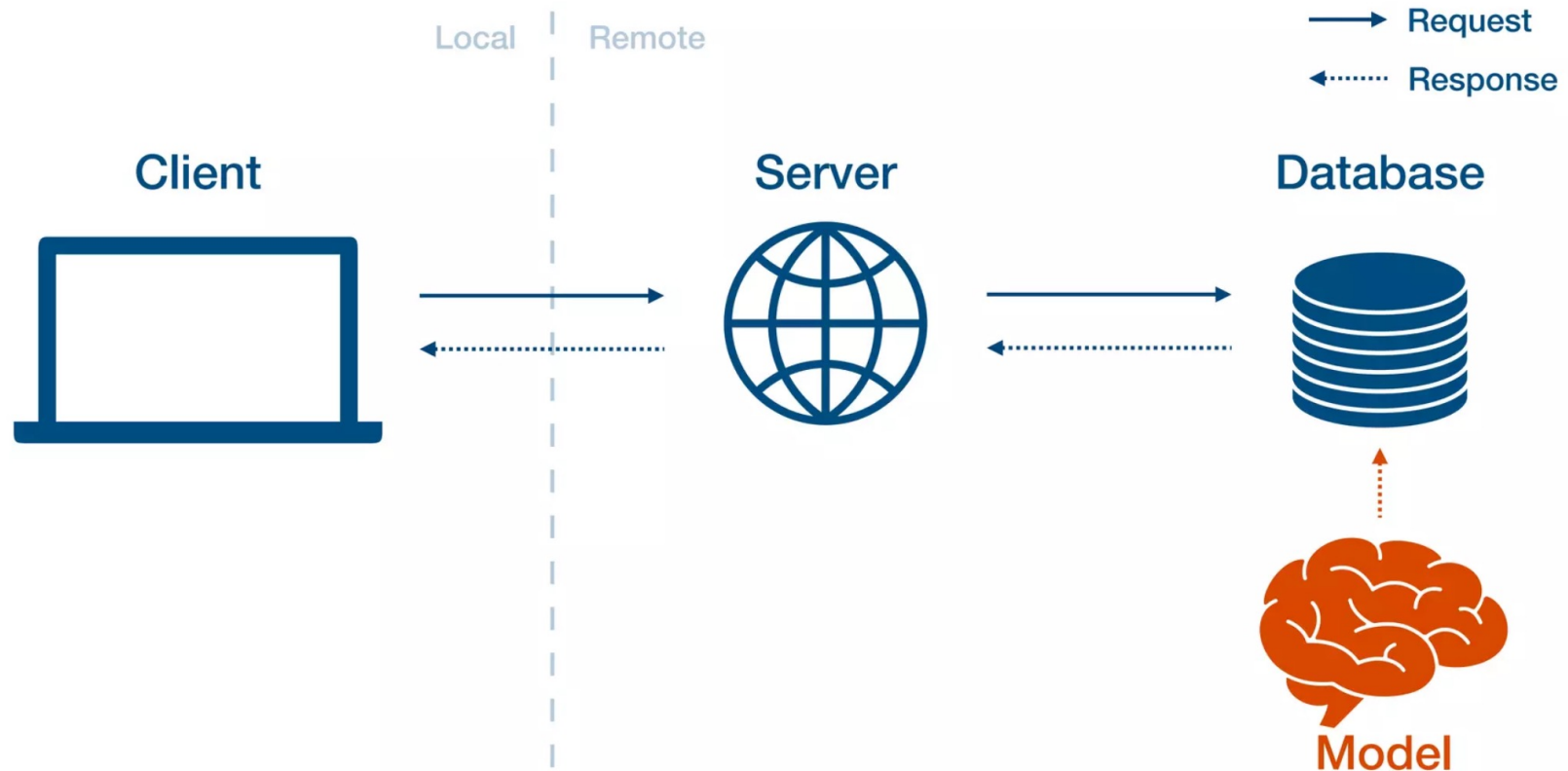
**You know, I'm something  
of a software engineer myself**

imgflip.com

# Where to host the model?

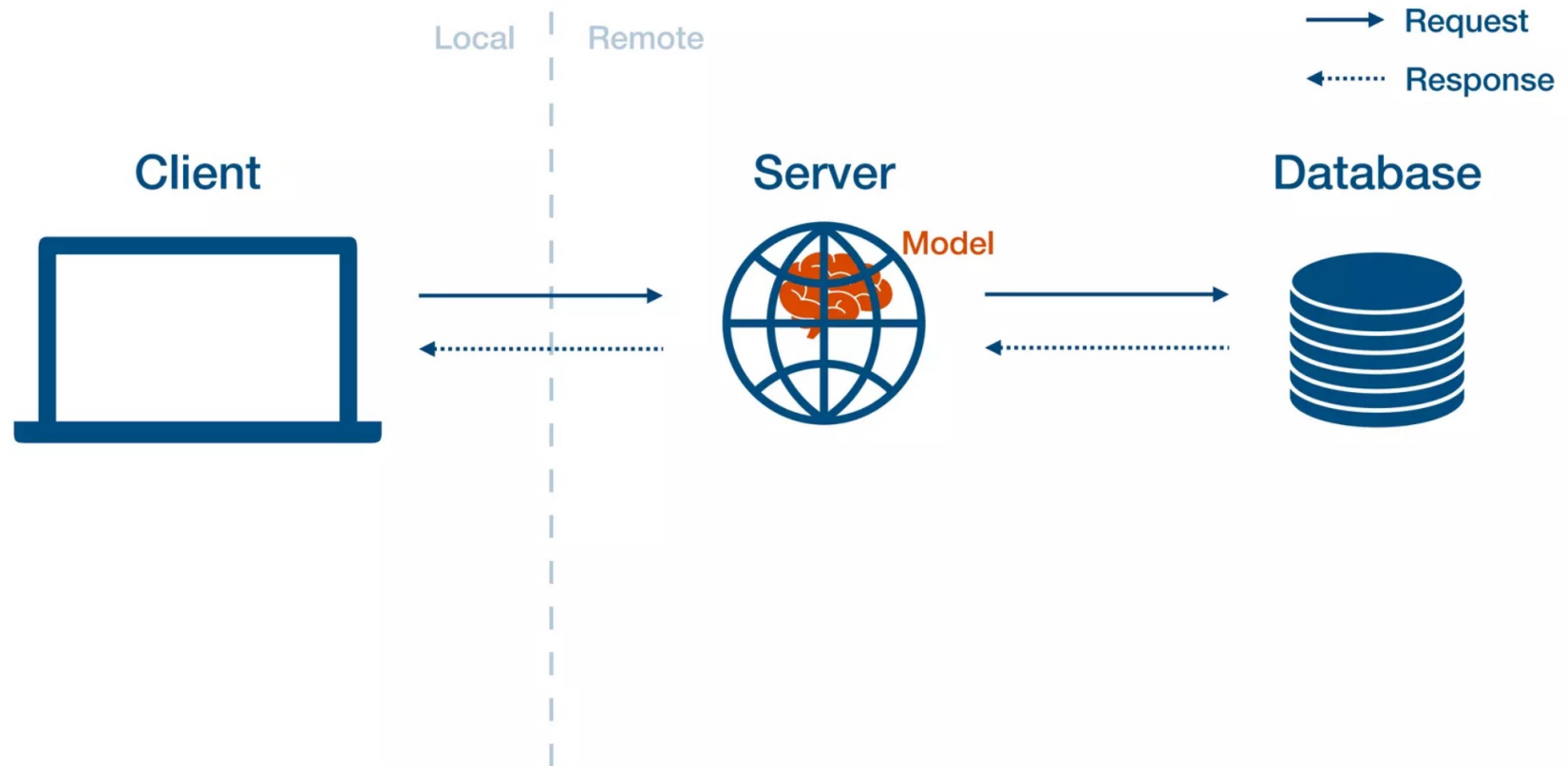


# Batch/offline prediction



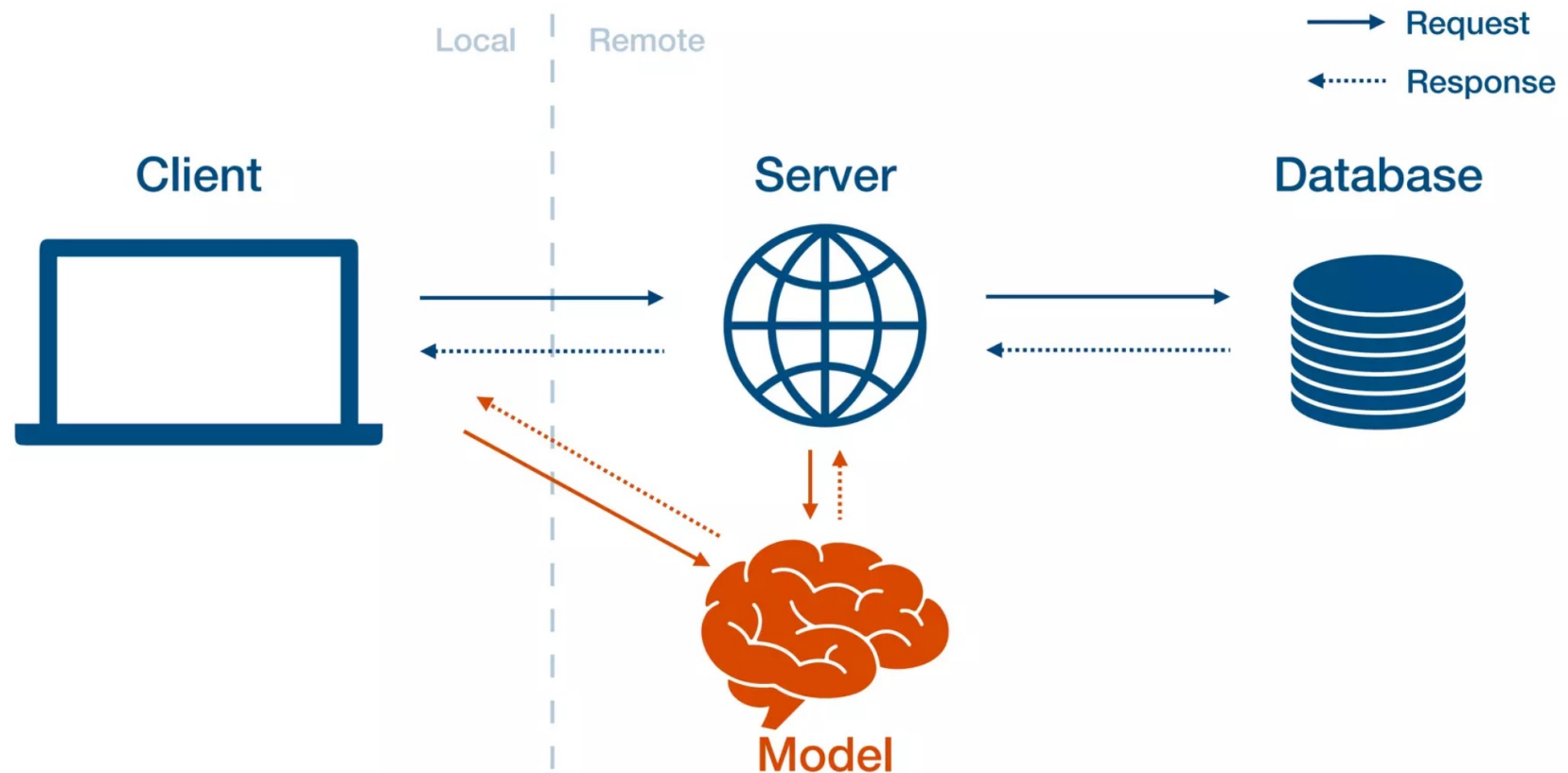
- Periodically run your model on new data and cache the results in a database
- Works if the universe of inputs is relatively small (e.g., 1 prediction per user)

# Realtime/Online prediction



Embedded within an application? What if the application is in Java?

# Model-as-a-Service



- Run your model on its own web server
- The backend (or the client itself) interact with the model by making requests to the model service and receiving responses back

# Want more info on deployment?

140 (excellent) slides with associated videos to be found here:  
<https://fullstackdeeplearning.com/spring2021/lecture-11/>



Week	Topic and Papers	Slides, Videos; Review Forms, Deadlines
0	<b>Introduction, ML Lifecycle Overview, and Basics</b>	Slides: <a href="#">PDF</a> <a href="#">PPTX</a> Video 1; Video 2; Video 3
	Readings: <a href="#">SIGMOD tutorial 1</a> , <a href="#">SIGMOD tutorial 2</a> , <a href="#">Berkeley report</a>	
1-2	<b>Topic 1: Classical ML Training at Scale</b>	Slides: <a href="#">PDF</a> <a href="#">PPTX</a> Video 1; Video 2
	For review: <a href="#">Parameter Server</a>	<a href="#">Review 1 Form</a> ; due 10/6
	For review: <a href="#">XGBoost</a>	<a href="#">Review 2 Form</a> ; due 10/13
	More readings: <a href="#">MADlib</a> , <a href="#">MLlib</a> , <a href="#">Mahout</a> , <a href="#">GraphLab</a> , <a href="#">AWS Sagemaker</a>	-
1	<b>No class on 10/8</b>	-
3	<b>Topic 2: Deep Learning Systems</b>	Slides: <a href="#">PDF</a> <a href="#">PPTX</a> Video 1; Video 2; Video 3
	For review: <a href="#">TensorFlow</a> (Talk slides)	<a href="#">Review 3 Form</a> ; due 10/20
	More readings: <a href="#">Horovod</a> , <a href="#">Distributed PyTorch</a> , <a href="#">TVM</a>	-
4-5	<b>Topic 3: Feature Engineering and Model Selection Systems</b>	Slides: <a href="#">PDF</a> <a href="#">PPTX</a> Video 1 Video 2
	For review: <a href="#">Cerebro</a>	<a href="#">Review 4 Form</a> ; due 10/27
	More readings: <a href="#">MSMS</a> , <a href="#">Hyperband</a> , <a href="#">ASHA</a> , <a href="#">Vizier</a> , <a href="#">Columbus</a> , <a href="#">Vista</a>	-
5	<b>Review Session 1 on 11/3 (tentative)</b>	Slides: <a href="#">PDF</a>
5	<b>Exam 1 on 11/5</b>	-
6	<b>Topic 4: Data Sourcing and Organization for ML</b>	Slides: <a href="#">PDF</a> <a href="#">PPTX</a> Video 1; Video 2; Video 3
	For review: <a href="#">TFDV</a>	<a href="#">Review 5 Form</a> ; due 11/3
	More readings: <a href="#">Deequ</a> , <a href="#">Snorkel</a> , <a href="#">Ground</a> , <a href="#">SortingHat</a> , <a href="#">Hamlet</a>	-
7	<b>Guest Lecture by Matei Zaharia (Databricks and Stanford) on MLFlow on 11/17</b>	<a href="#">Video</a> ; <a href="#">Slides</a> <a href="#">PDF</a>
7-9	<b>Topic 5: ML Deployment</b>	Slides: <a href="#">PDF</a> <a href="#">PPTX</a> Video 1; Video 2
	For review: <a href="#">Clipper</a>	<a href="#">Review 6 Form</a> ; due 11/12
	More readings: <a href="#">TF Serving</a> , <a href="#">Uber PyML</a> , <a href="#">Hummingbird</a> , <a href="#">Federated ML</a>	-
8	<b>Guest Lecture by Angela Jiang (Determined AI) on Determined DL Platform on 11/24</b>	<a href="#">Video</a> ; <a href="#">Slides</a> <a href="#">PDF</a>
8	<b>Thanksgiving Holiday on 11/26</b>	-
9	<b>Guest Lecture by Joshua Patterson (NVIDIA) on RAPIDS on 12/1</b>	<a href="#">Video</a> ; <a href="#">Slides</a> <a href="#">PDF</a>
9-10	<b>Topic 6: ML Platforms and Feature Stores</b>	Slides: <a href="#">PDF</a> <a href="#">PPTX</a> Video1; Video 2
	For review: <a href="#">ML systems technical debt</a>	<a href="#">Review 7 Form</a> ; due 11/17
	For review: <a href="#">TensorFlow Extended</a>	<a href="#">Review 8 Form</a> ; due 12/3
	More readings: <a href="#">MLFlow</a> , <a href="#">Michelangelo</a>	-

**CSE 234/291 from Fall'20 with lecture videos on Youtube**  
<https://cseweb.ucsd.edu/classes/fa20/cse291-d/schedule.html>



## DSC 102 focuses on thinking about the fundamentals of scalable analytics systems

1. **“Systems”**: What resources does a computer have?  
How to store and efficiently compute over large data?  
What is cloud?
2. **“Scalability”**: How to scale and parallelize data-intensive computations?
3. **For “Analytics”**:
  1. **Source**: Data acquisition & preparation for ML
  2. **Build**: Model selection & deep learning systems
  3. **Deploying** ML models
4. Hands-on experience with scalable analytics tools