

Fine-Grained Long-Range Prediction of Resource Usage

Ivan, Petru, and Zebo

March 15, 2017

Outline

1. Goal

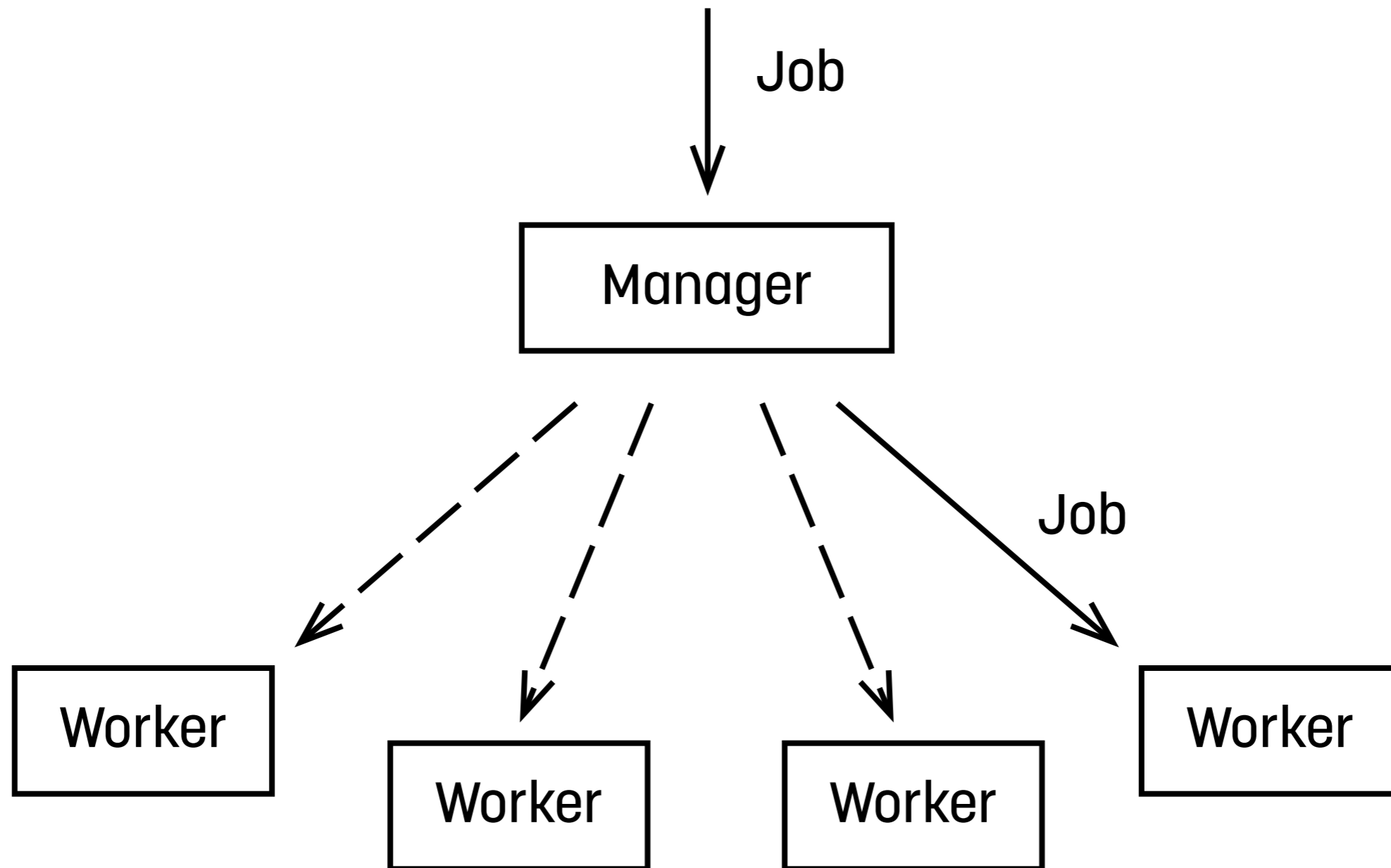
2. Data

3. Model

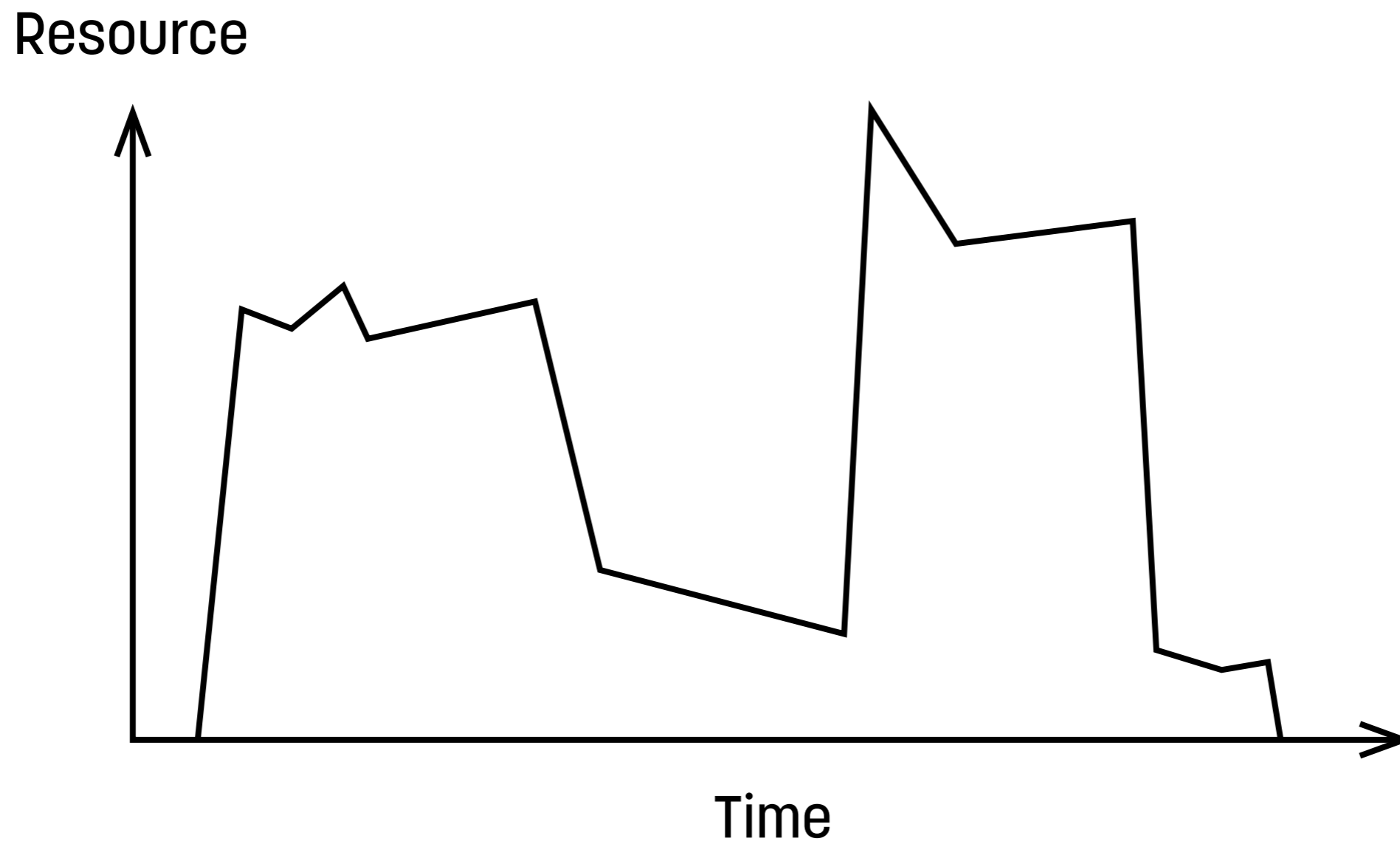
4. Tuning

Goal

Scenario



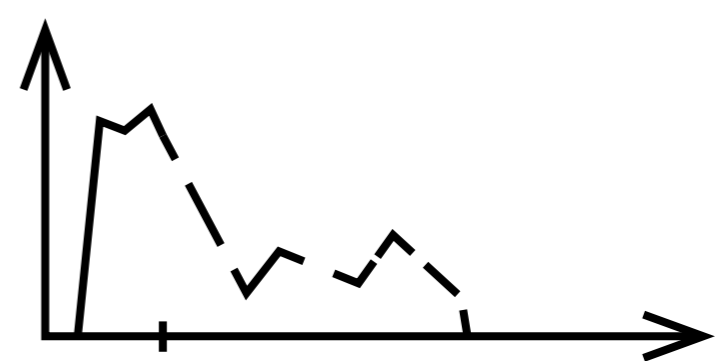
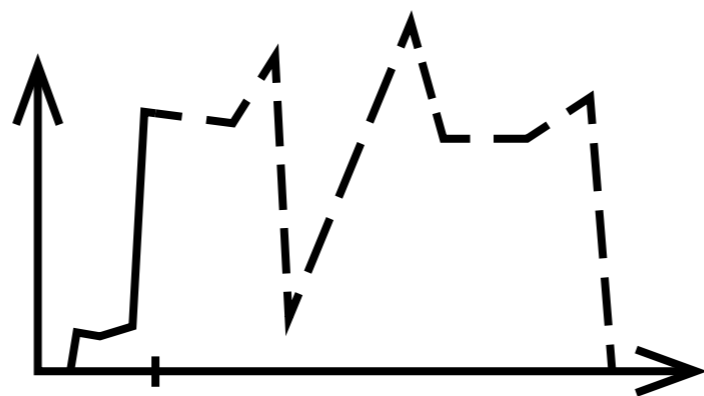
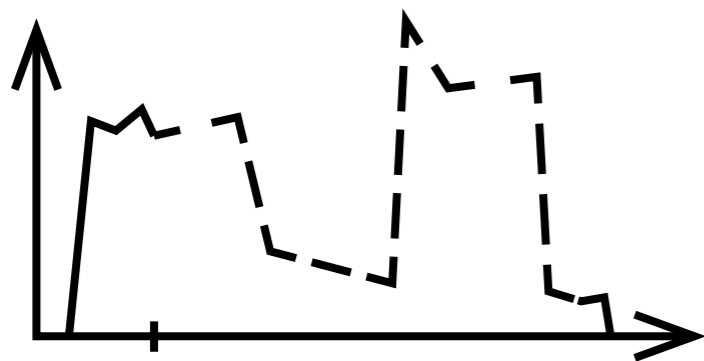
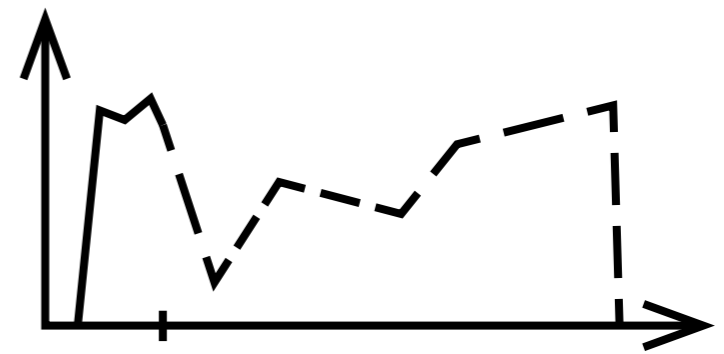
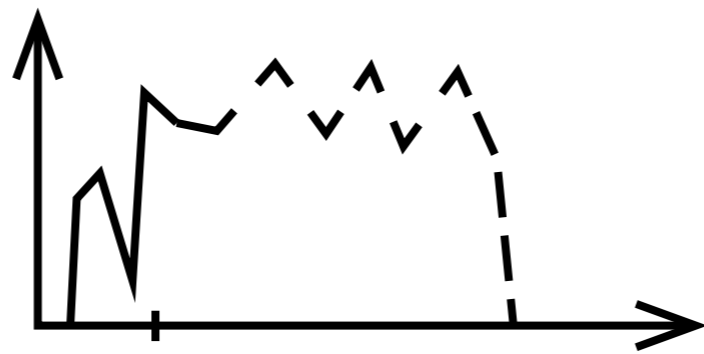
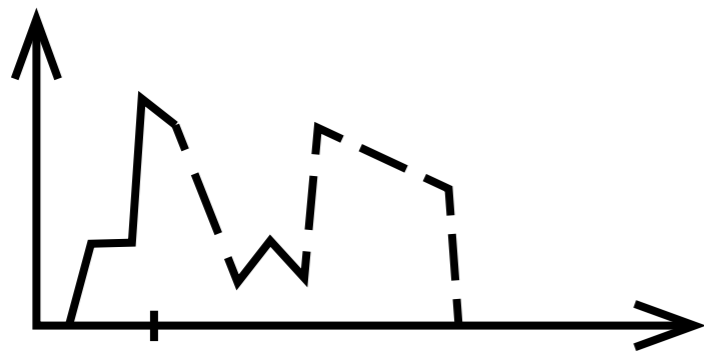
Scenario



Premise

- Knowing the future is useful

Prediction



Objective

- Predict the future resource usage

Means

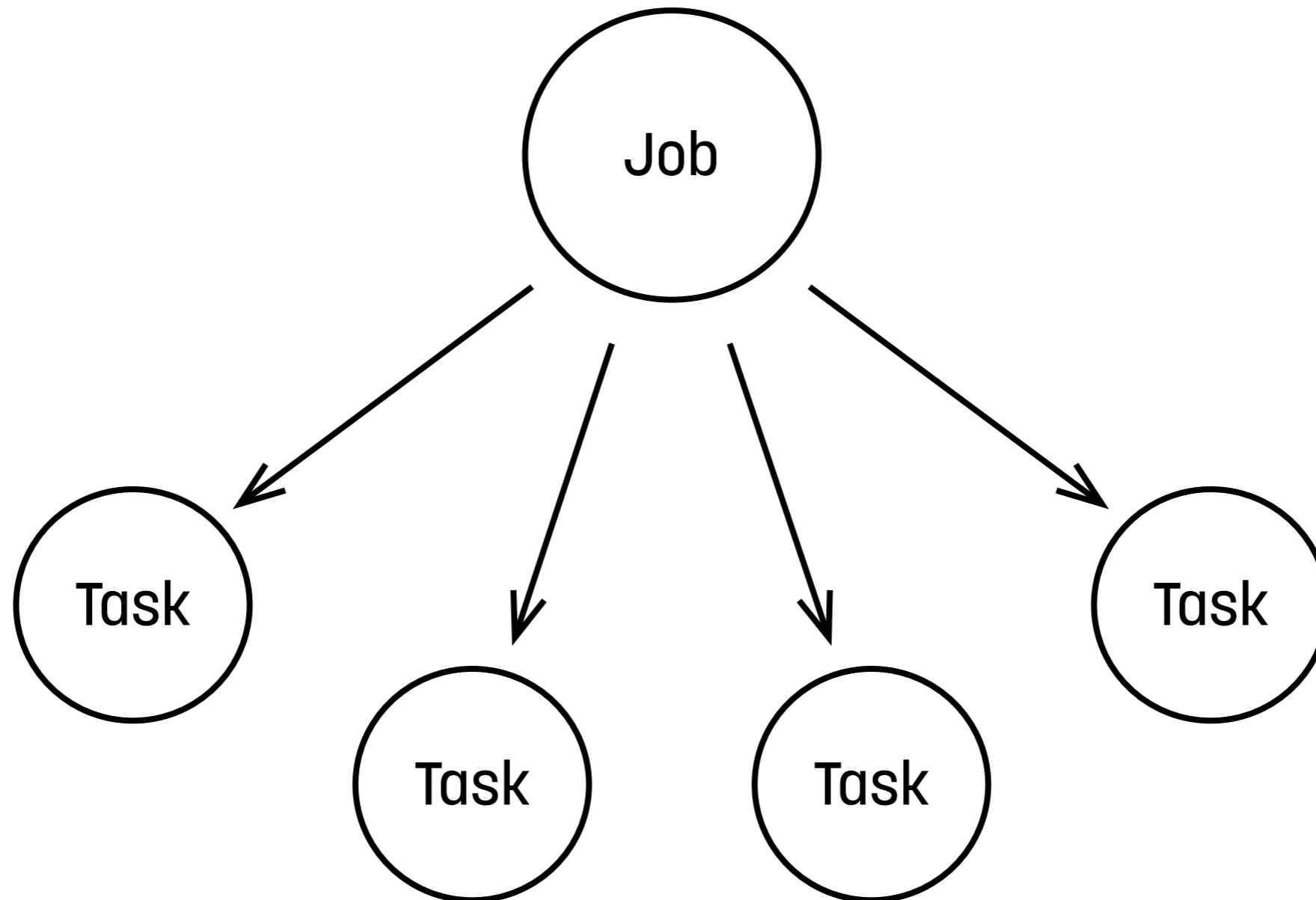
- Machine-learn from data

Data

Google Cluster Usage

- 1 month
- 1,000 users
- 700,000 jobs
- 13,000 machines

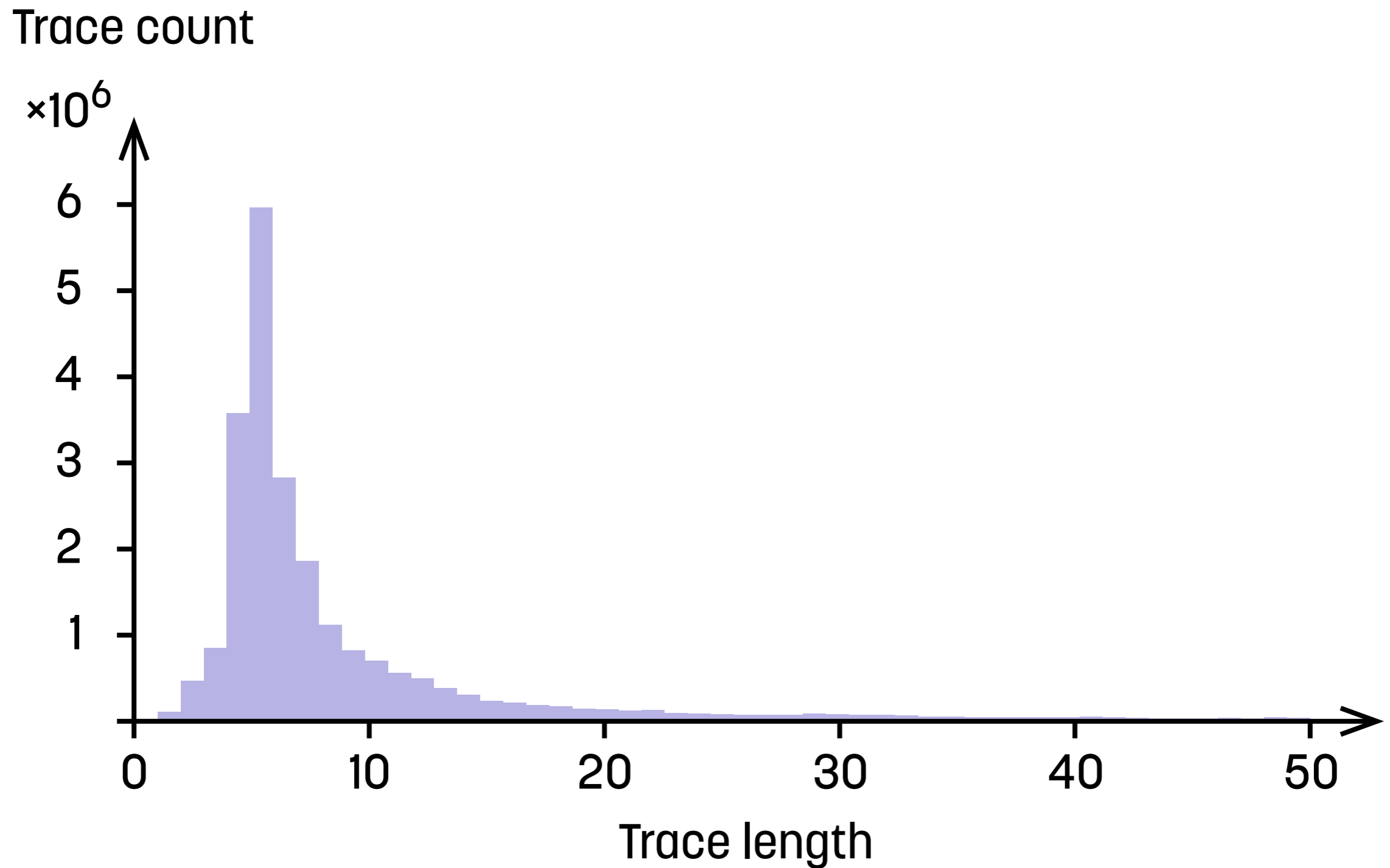
Google Cluster Usage



Task Resource Usage

- 200 GB
- 25,000,000 tasks
- 1,300,000,000 records
- CPU, memory, and disk usage
- Max and mean over 5-minute intervals

Task Resource Usage



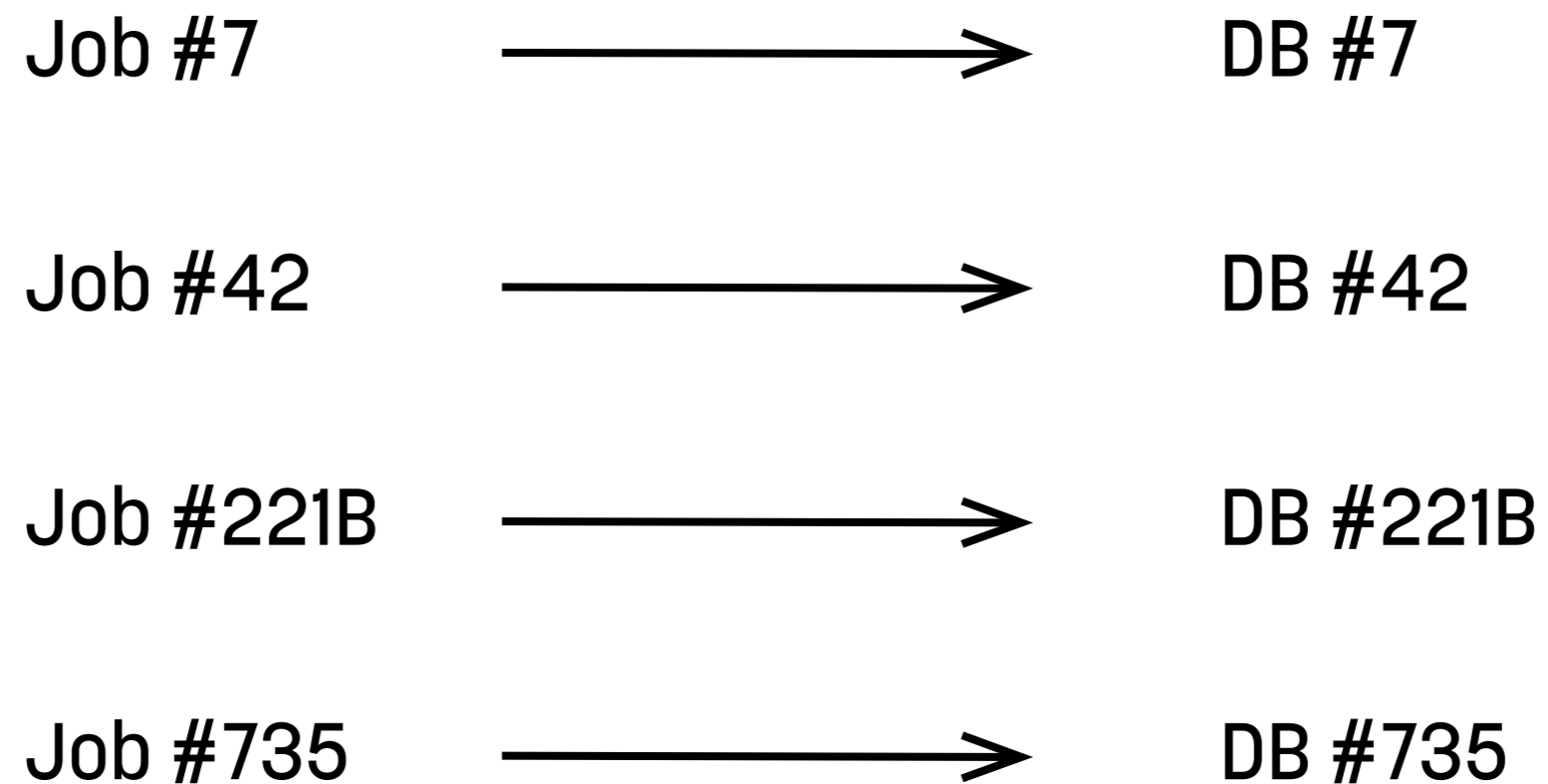
Problem

- Need fast access to individual traces to streamline machine learning

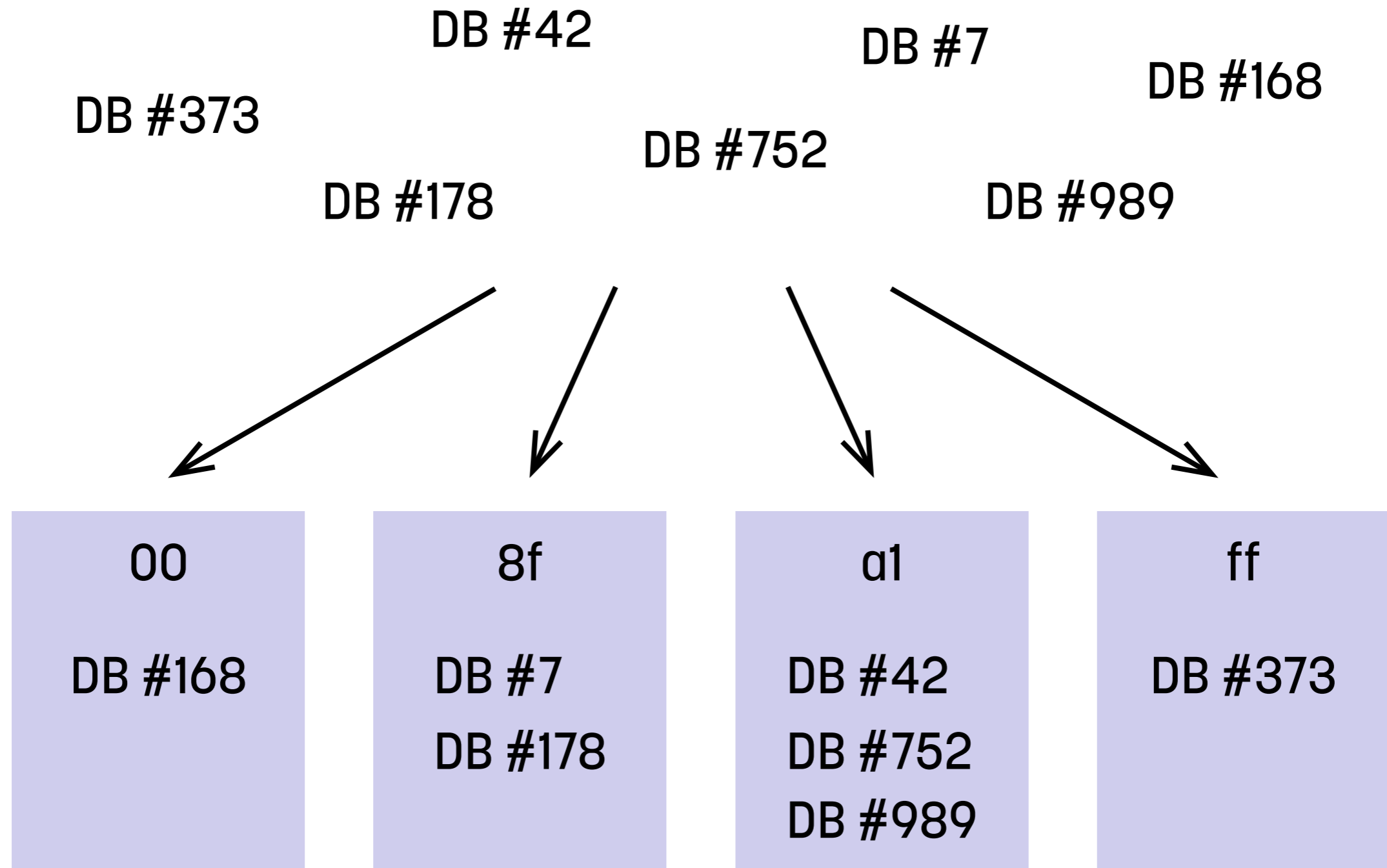
Divide & Conquer

- Preselect individual traces and store them in separate databases

Divide & Conquer



Divide & Conquer



Model

Problem

- Design an adequate predictive model

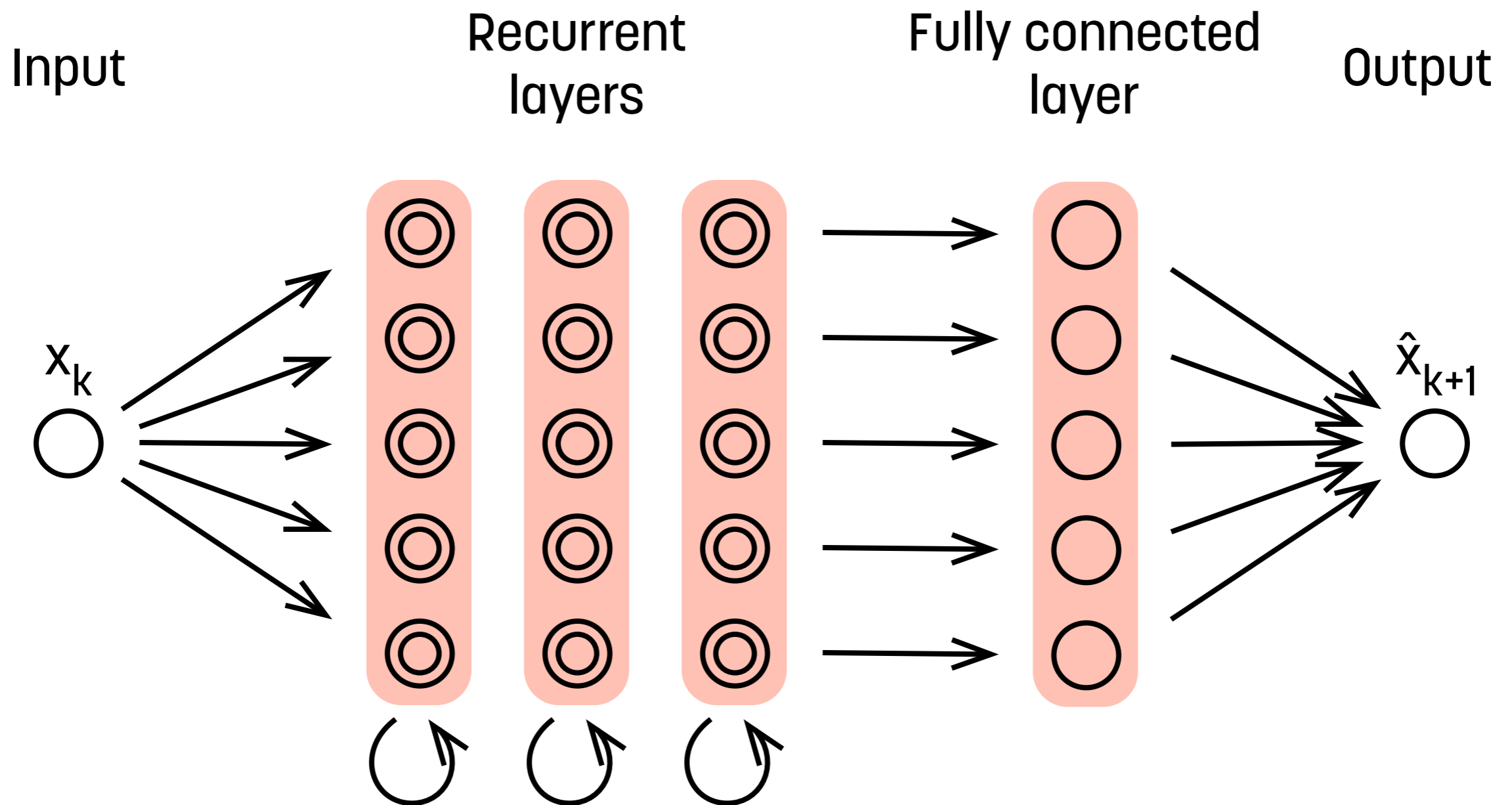
Premise

- Artificial neural networks are the state-of-the-art

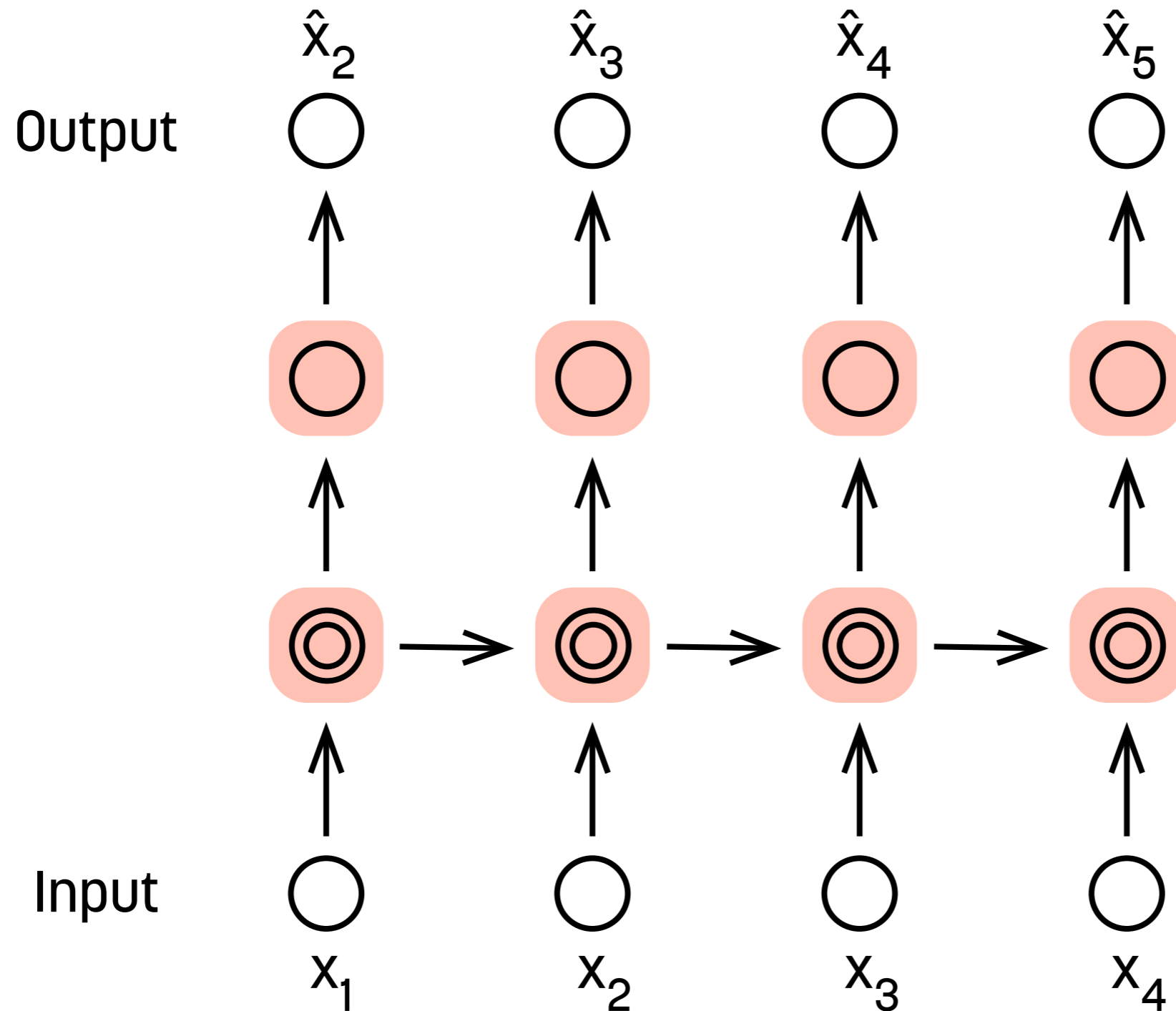
Architecture

- [The Neural Network Zoo \(click me\)](#)

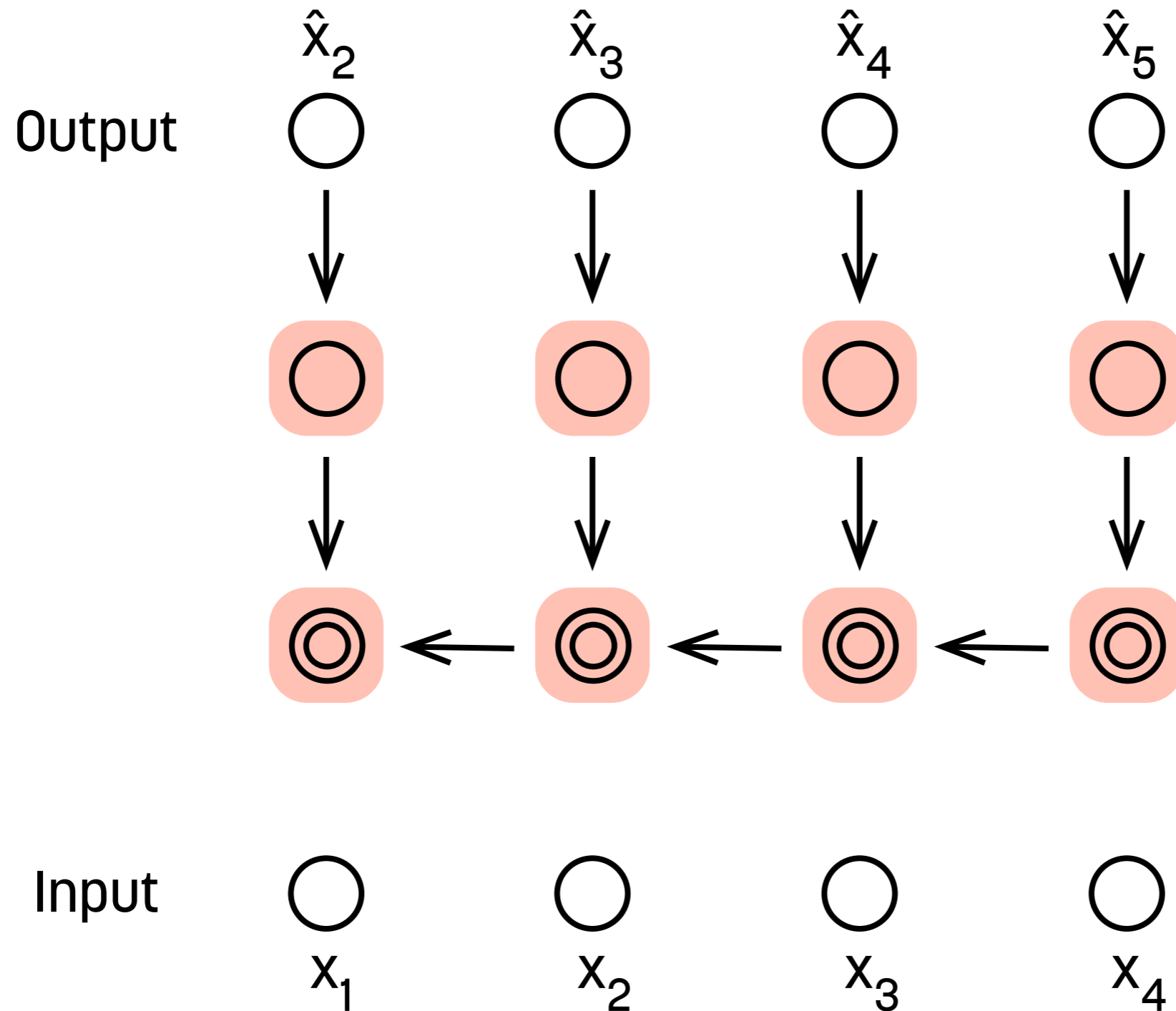
Architecture



Unrolling



Training



Software

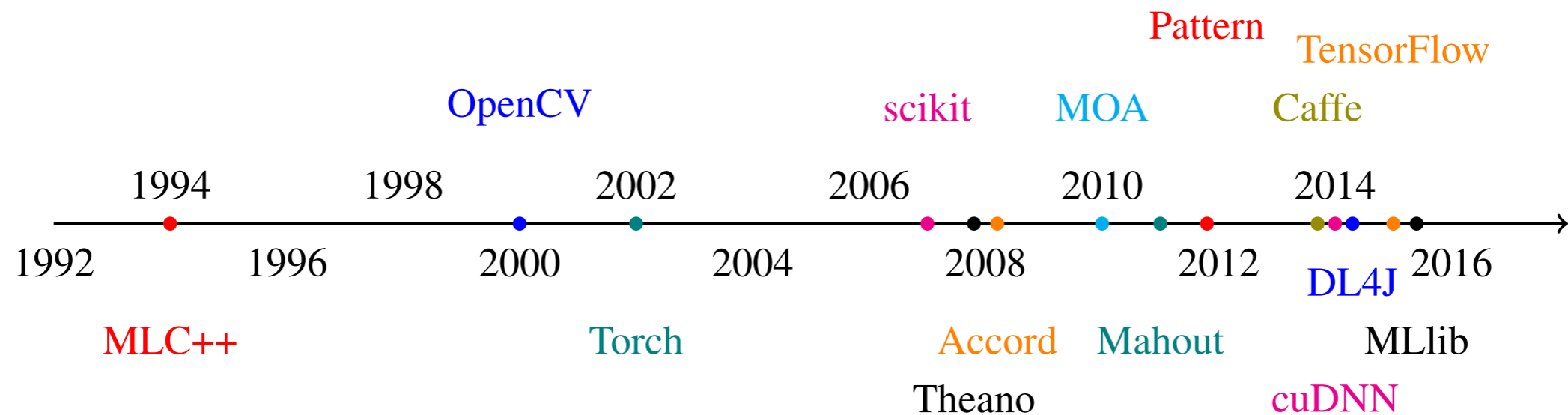
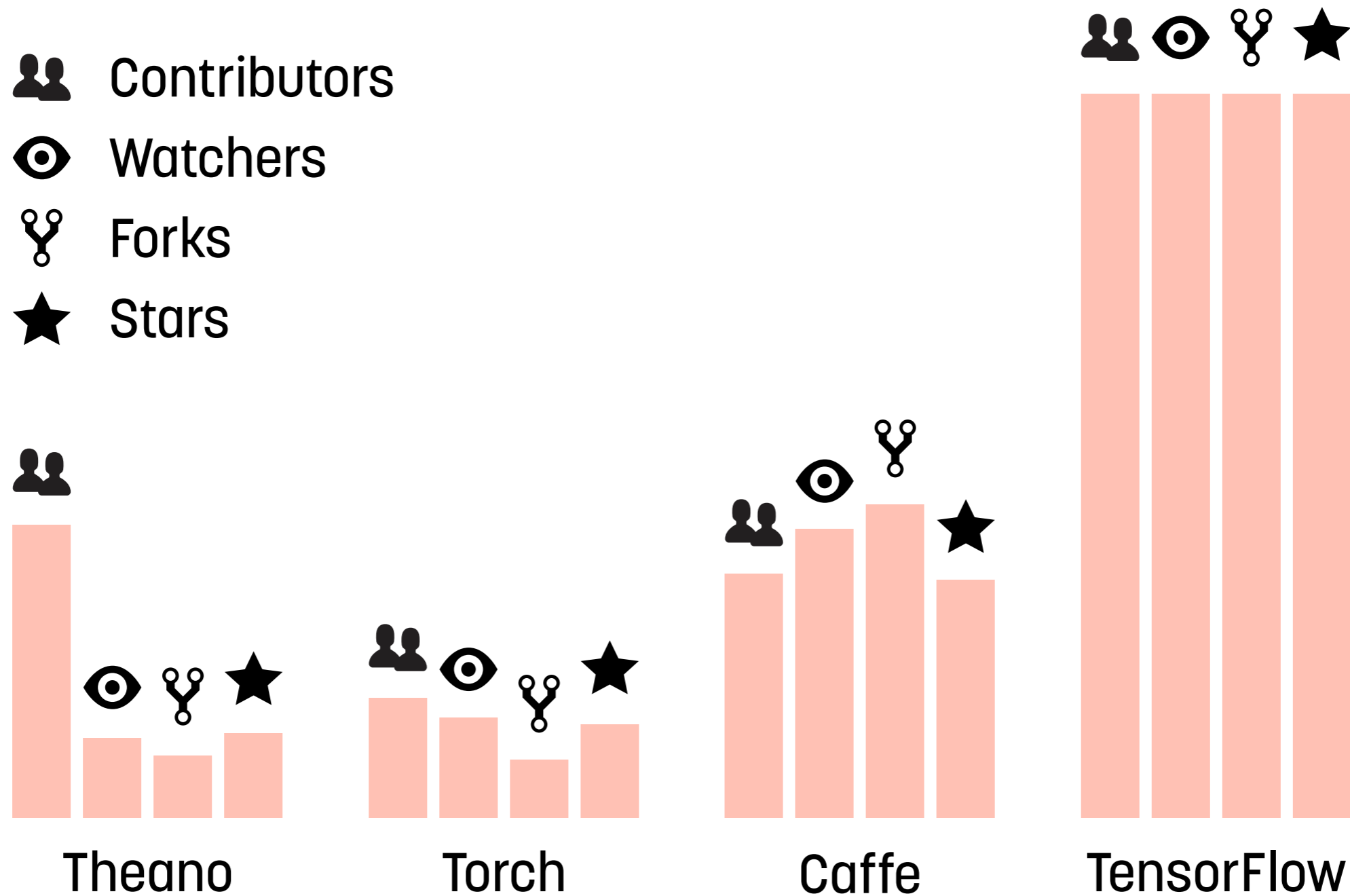


Fig. 1: A timeline showing the release of machine-learning libraries discussed in section I in the last 25 years.

Software



Source: <https://github.com>

TensorFlow

- Google Brain
- Open source
- Flexible & efficient
- Scalable & portable
- User-friendly (show demo)

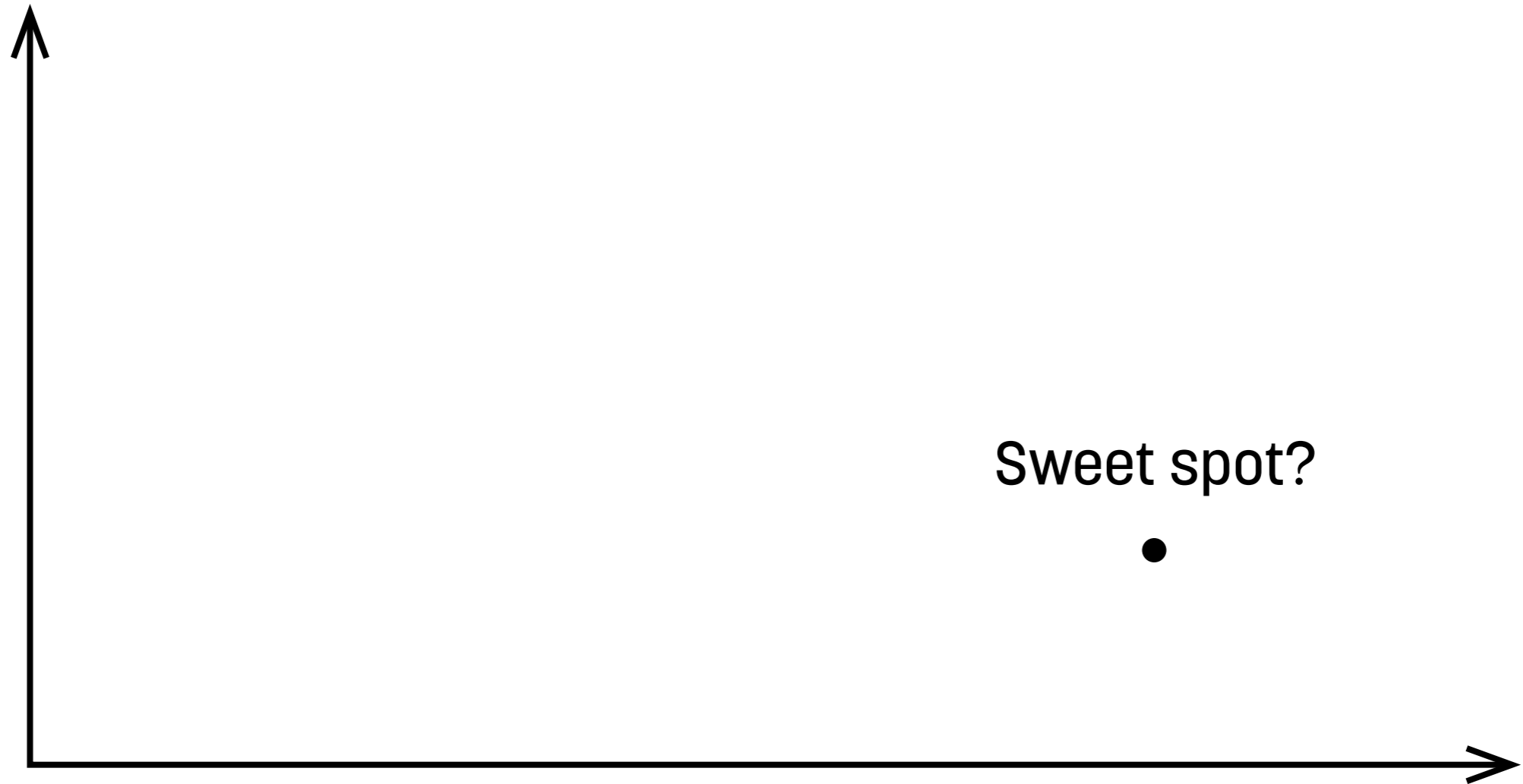
Tuning

Problem

- Decide on the hyperparameters

Configuration Selection

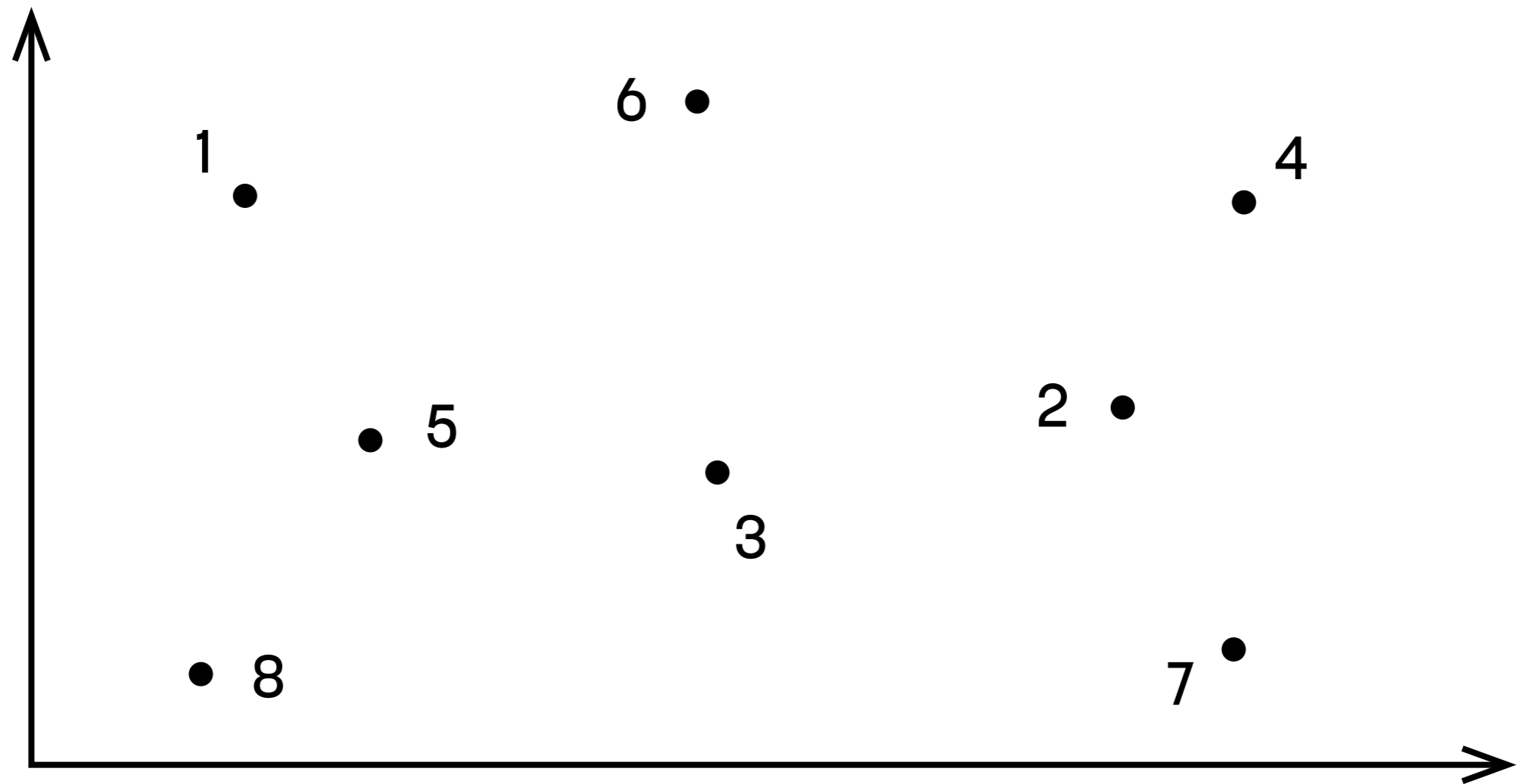
Parameter 2



Parameter 1

Exploration

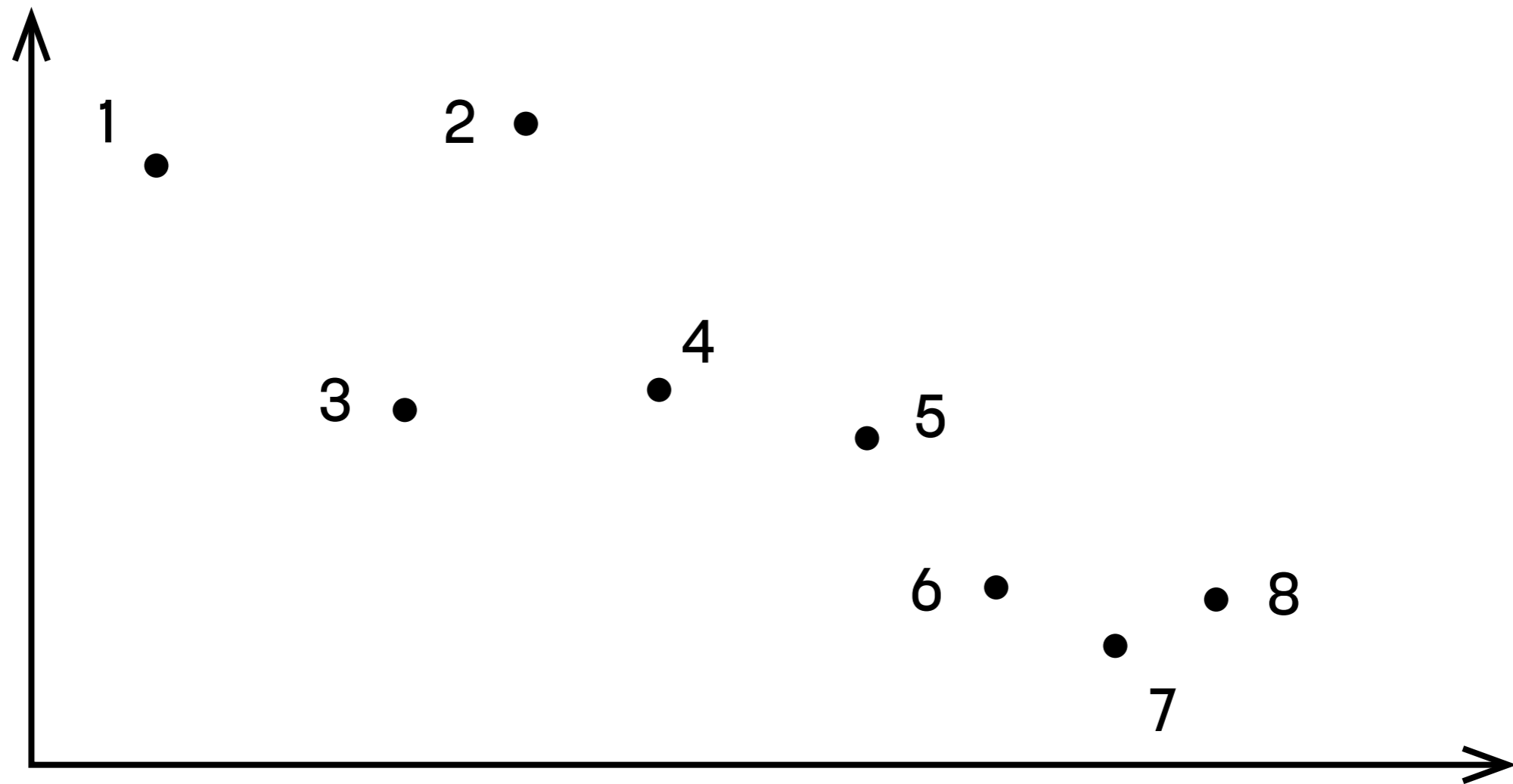
Parameter 2



Parameter 1

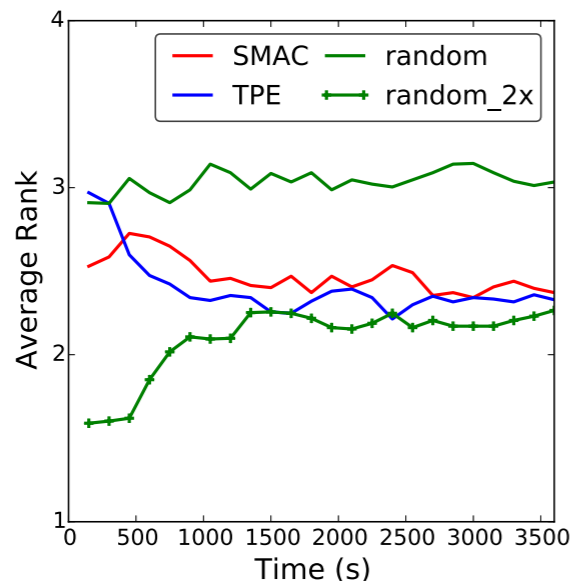
Exploitation

Parameter 2

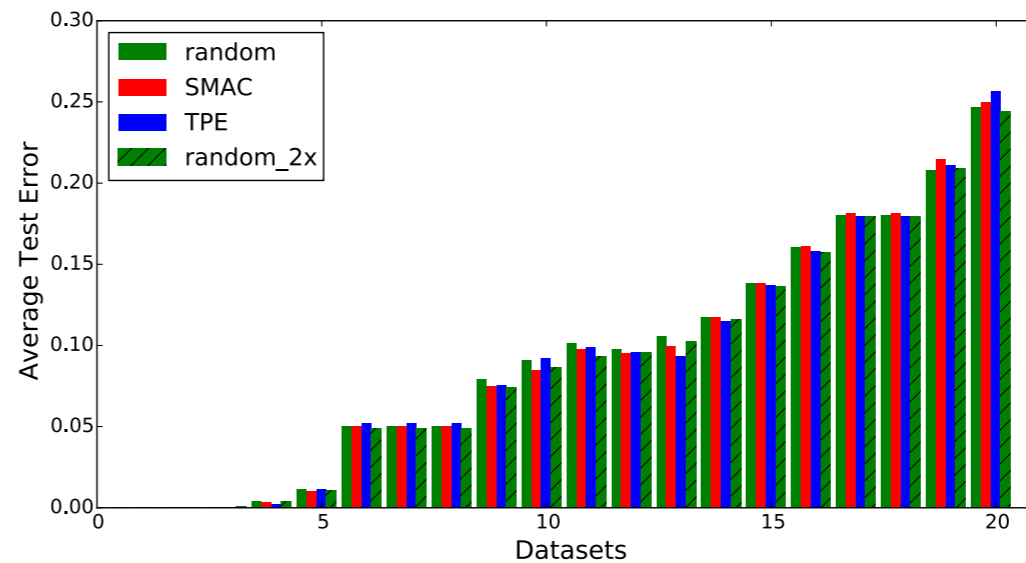


Parameter 1

Exploration & Exploitation



(a) Rank for 117 Datasets



(b) Test Error for 20 Datasets

Figure 2: Empirical evaluation of various search methods on 117 datasets. Search methods were executed for a one hour duration for each dataset, continuously reporting their best identified models throughout this time window. Models were evaluated using an unseen test set. Results are reported for random search ('random'), random search run on two machines ('random_2x'), and two Bayesian optimization methods ('SMAC', 'TPE'). (a) Average rank of test error across all datasets, where lower is better. The rank for each dataset is based on the average test error across 20 trials. (b) Average test error for 20 randomly sampled datasets after one hour of execution. See Figure A.3 for corresponding results for all 117 datasets.

Hyperband

- Pure exploration
- Adaptive resource allocation
- 5-30× speedup over state-of-the-art Bayesian optimization algorithms

Hyperband

i	$s = 4$		$s = 3$		$s = 2$		$s = 1$		$s = 0$	
	n_i	r_i	n_i	r_i	n_i	r_i	n_i	r_i	n_i	r_i
0	81	1	27	3	9	9	6	27	5	81
1	27	3	9	9	3	27	2	81		
2	9	9	3	27	1	81				
3	3	27	1	81						
4	1	81								

Table 1: The values of n_i and r_i for the brackets of HYPERBAND corresponding to various values of s , when $R = 81$ and $\eta = 3$.

Conclusion

1. Goal

2. Data

3. Model

4. Tuning

Thank you!
Questions?