

# Systems plus ML: When the sum is greater than its parts

Ion Stoica

UC Berkeley, Director of RISELab

November 5, 2019





Studies the design of  
**r**real-time,  
**i**ntelligent,  
**s**ecure, and  
**e**xplainable  
algorithms and systems





Studies the design of

**r**real-time,

**i**ntelligent,

**s**ecure, and

**e**xplainable

algorithms and systems

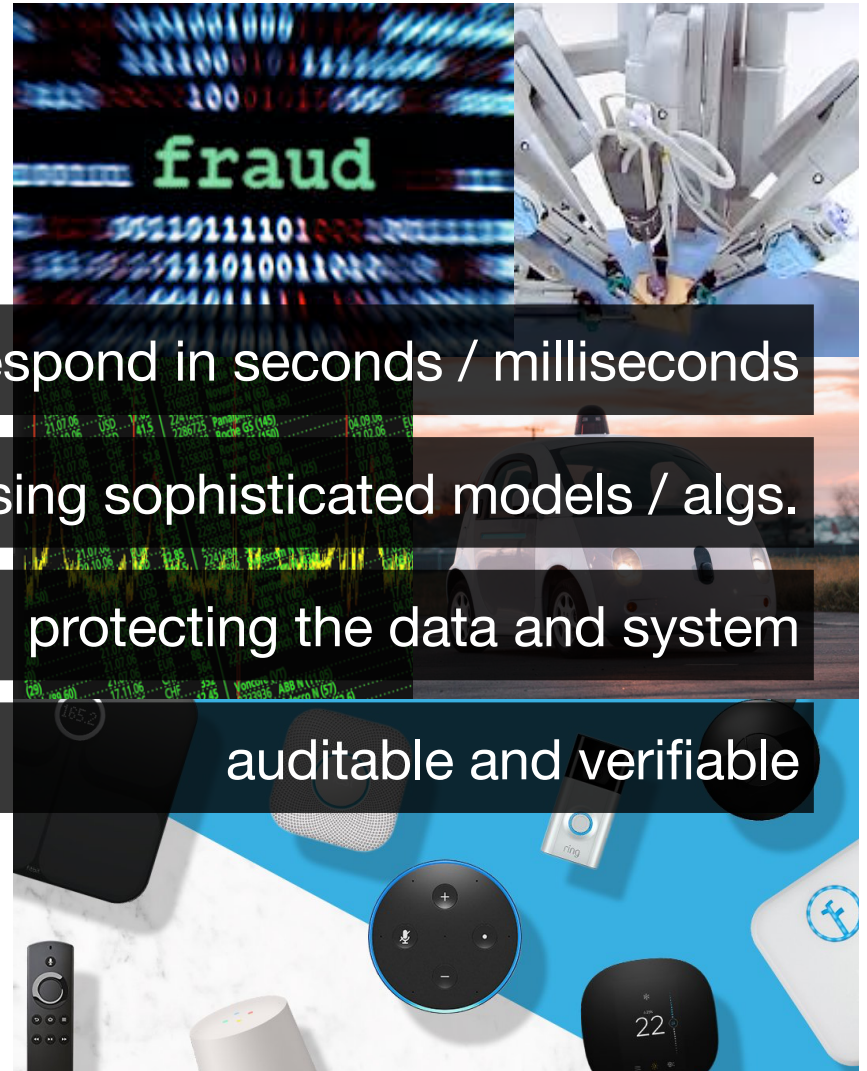


respond in seconds / milliseconds

using sophisticated models / algs.

protecting the data and system

auditable and verifiable





Studies the design of

**r**real-time,

**i**ntelligent,

**s**ecure, and

**e**xplainable

algorithms and systems

## Interdisciplinary Lab

AI



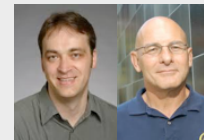
Abbeel Goldberg Gonzalez Jordan Mahoney

Security



Joseph Popa Song

Hardware



Asanovic Patterson

Systems



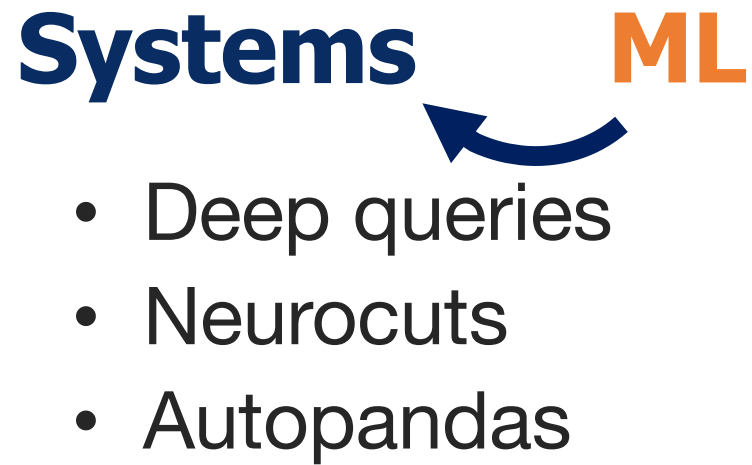
Ghodsi Culler Hellerstein Katz Stoica

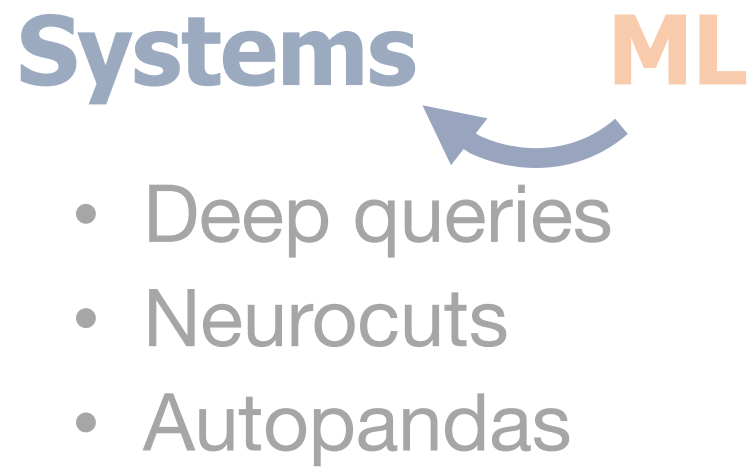
Collaborate with:





Positive feedback loop to accelerate ML and Systems





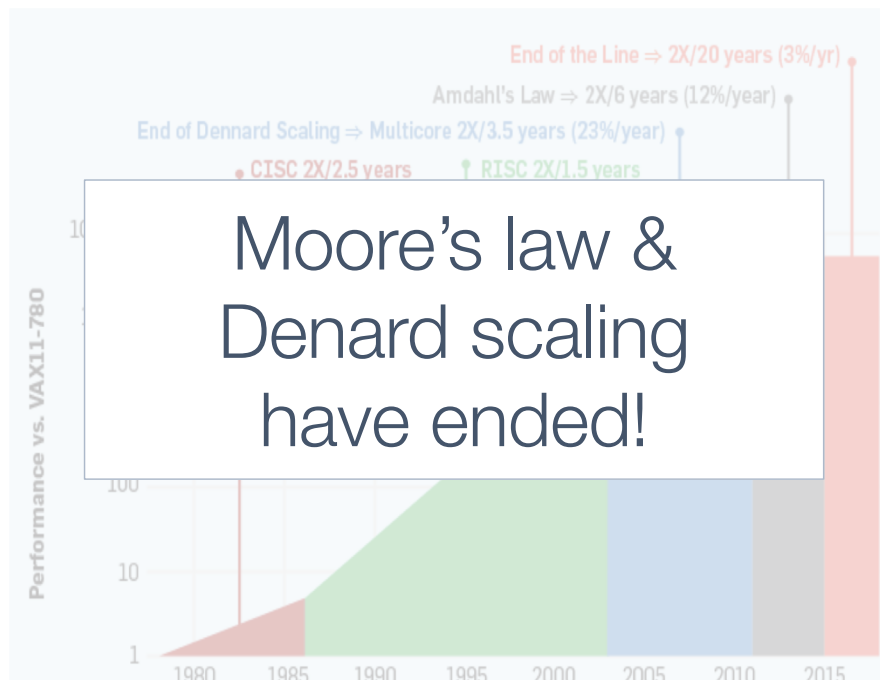
# Trends

Apps becoming **distributed**

Apps becoming more **complex**



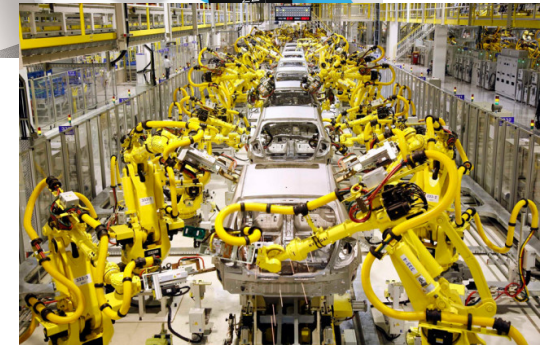
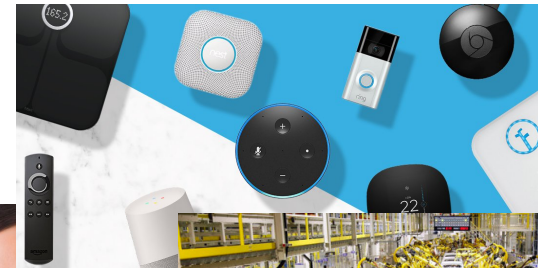
# Apps becoming distributed



No choice but to distribute apps

# Apps becoming more complex

Virtually all apps will become AI centric



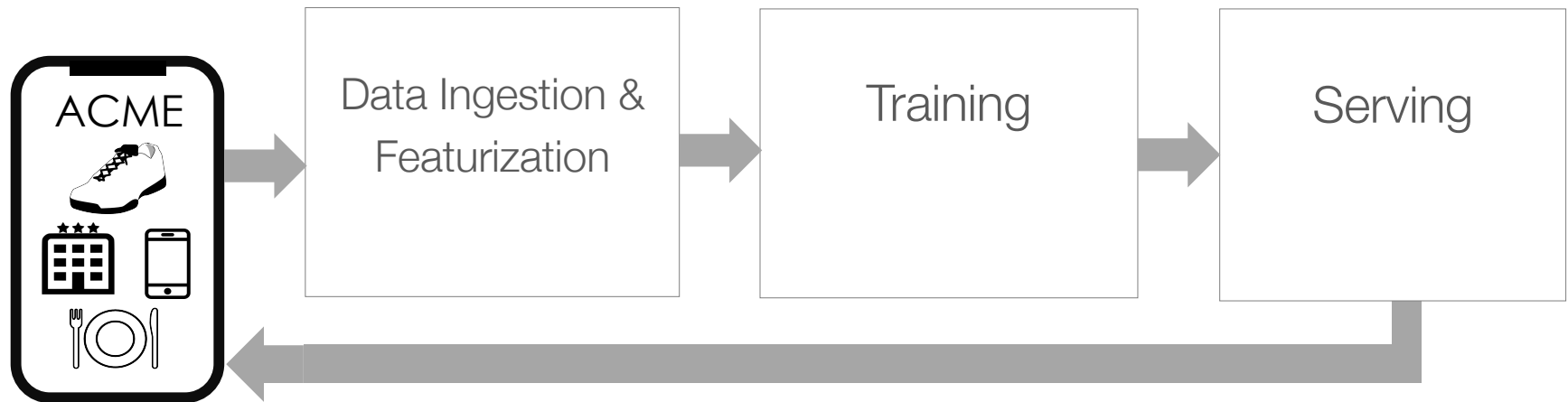
# Example: In-app promotion



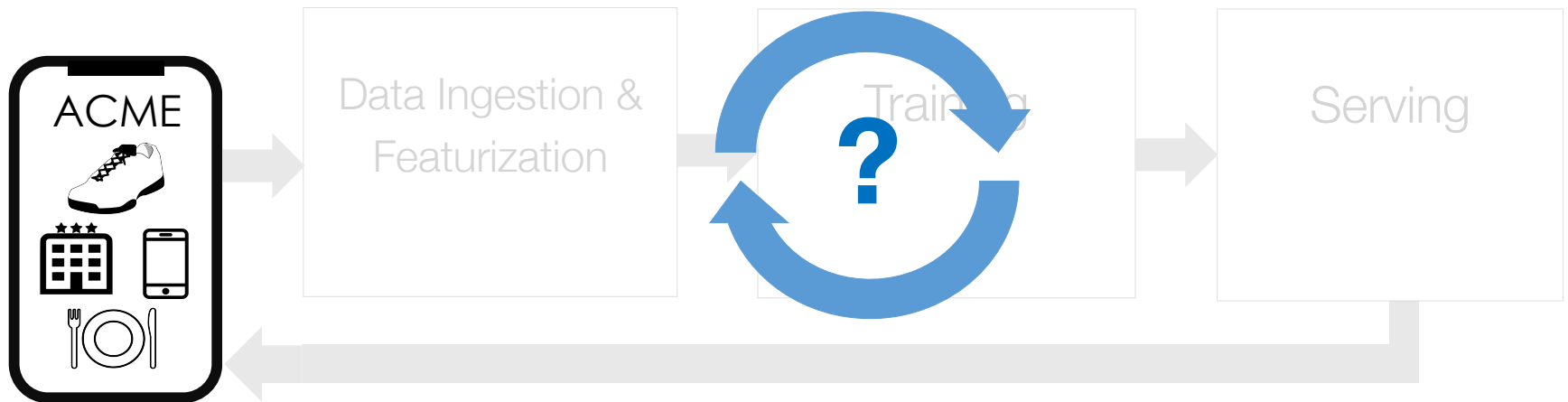
A real use case:

- Recommend services, products
- Largest fintech company in the world

# Example: In-app promotion



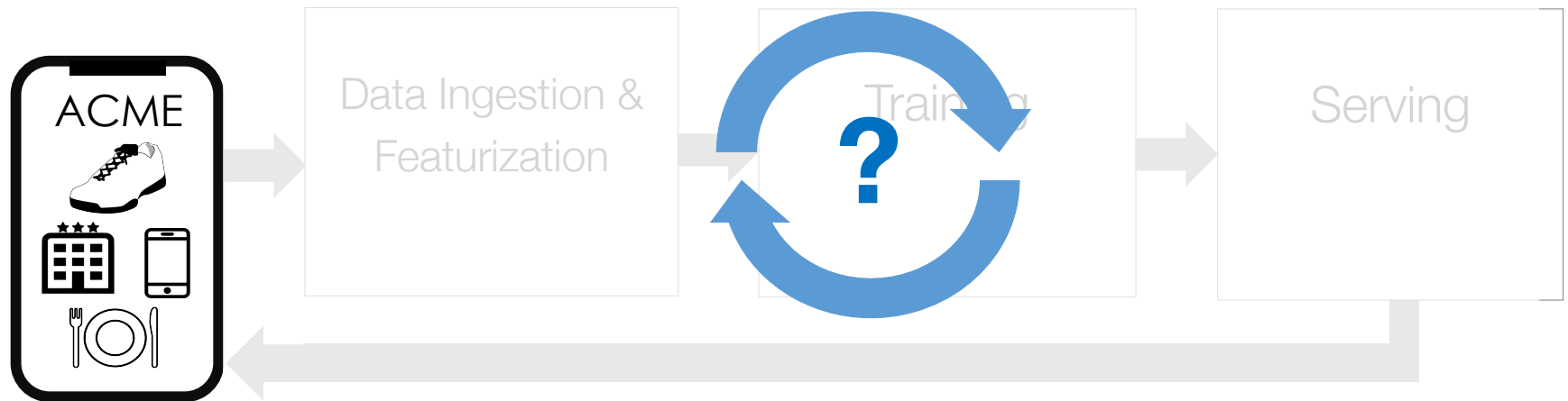
# Example: In-app promotion



Two questions:

- How fast can we update the model?
- How much does it matter?

# Example: In-app promotion



Model updated every 1 day



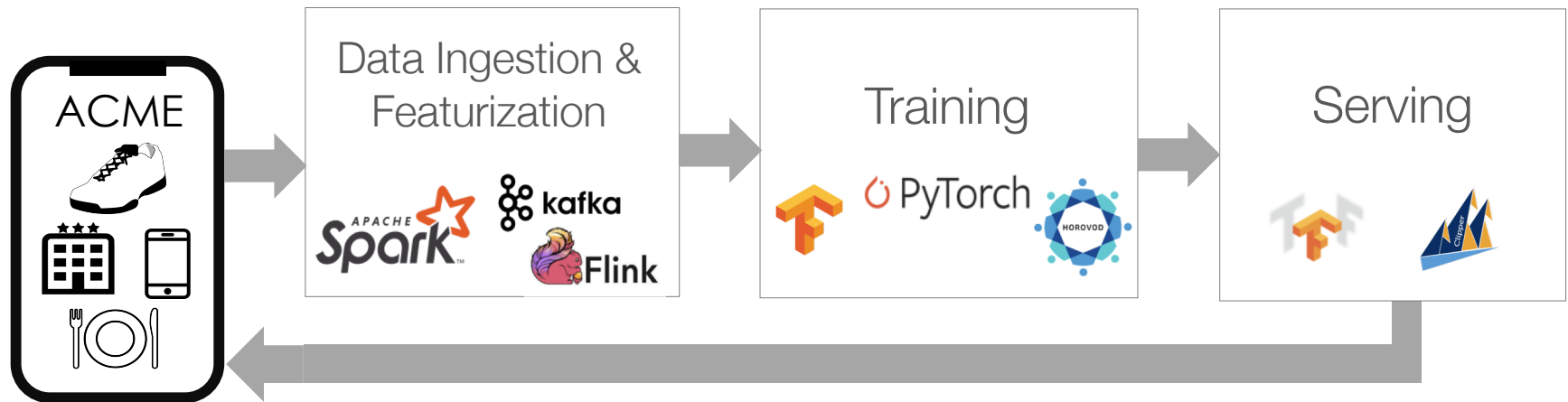
**+ 5%** CTR (Click Through Rate)

Model updated every 1 hour (state-of-the-art solution)



**Want to get lower, but how?**

# Example: In-app promotion



Previous solution: integrate best-of-breed frameworks

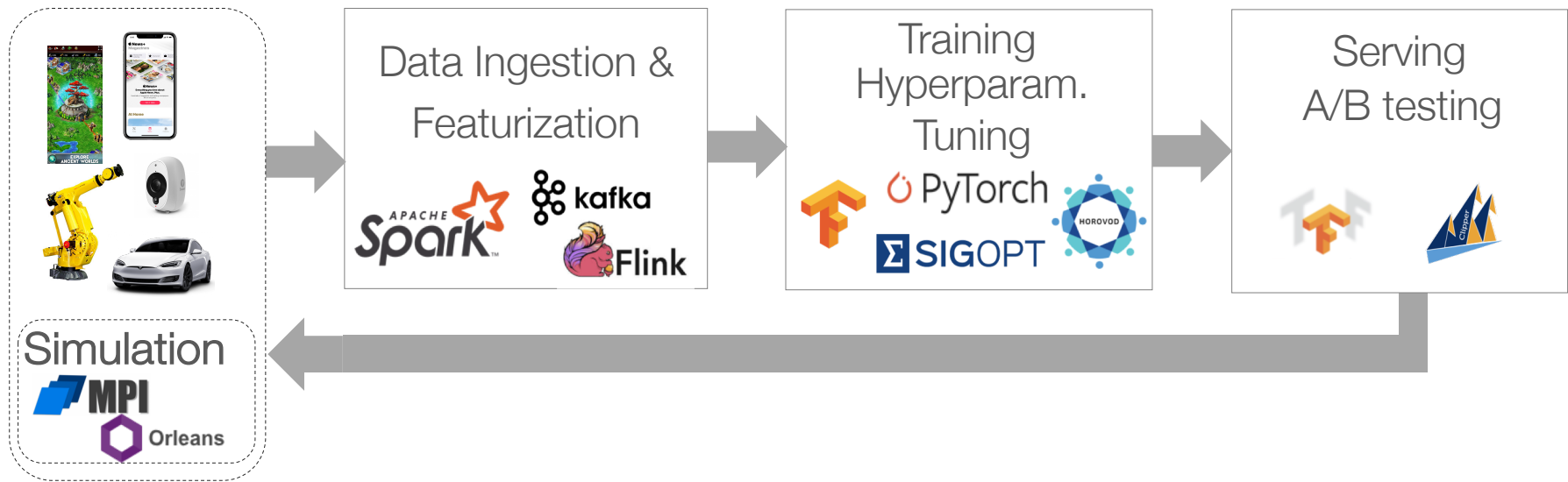
Challenges: end-to-end delay, development, management cost

# Even more complex patterns!





# Even more complex patterns!



## Reinforcement Learning

# Today's ML Ecosystem



**Distributed systems**

Training



**Distributed systems**

Model Serving



**Distributed systems**

Hyperparam. Tuning



**Distributed systems**

Streaming



**Distributed systems**

Simulation



**Distributed systems**

Featurization





## Libraries

Training

Model  
Serving

Hyperparam.  
Tuning

Streaming

Simulation

Featurization



# RAY

General-purpose distributed computing  
framework for Python (and Java)

\* "Ray: A Distributed Framework for Emerging AI Applications", Philipp Moritz et al, OSDI 2018

## Function

```
def read_array(file):  
    # read ndarray "a"  
    # from "file"  
    return a  
  
def add(a, b):  
    return np.add(a, b)  
  
a = read_array(file1)  
b = read_array(file2)  
sum = add(a, b)
```

## Object

```
class Counter(object):  
    def __init__(self):  
        self.value = 0  
    def inc(self):  
        self.value += 1  
        return self.value  
  
c = Counter()  
c.inc()  
c.inc()
```



## Function → Task

```
@ray.remote
def read_array(file):
    # read ndarray "a"
    # from "file"
    return a
```

```
@ray.remote
def add(a, b):
    return np.add(a, b)
```

```
a = read_array(file1)
b = read_array(file2)
sum = add(a, b)
```

## Object → Actor

```
@ray.remote
class Counter(object):
    def __init__(self):
        self.value = 0
    def inc(self):
        self.value += 1
        return self.value
```

```
c = Counter()
c.inc()
c.inc()
```

## Function → Task

```
@ray.remote
def read_array(file):
    # read ndarray "a"
    # from "file"
    return a

@ray.remote
def add(a, b):
    return np.add(a, b)

id1 = read_array.remote(file1)
id2 = read_array.remote(file2)
id = add.remote(id1, id2)
sum = ray.get(id)
```

## Object → Actor

```
@ray.remote
class Counter(object):
    def __init__(self):
        self.value = 0
    def inc(self):
        self.value += 1
        return self.value

c = Counter.remote()
id4 = c.inc.remote()
id5 = c.inc.remote()
```

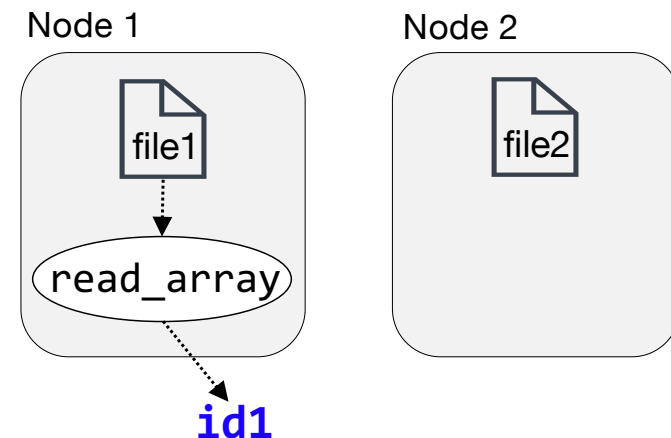
## Task API

```
@ray.remote
def read_array(file):
    # read ndarray "a"
    # from "file"
    return a
```

```
@ray.remote
def add(a, b):
    return np.add(a, b)
```

```
id1 = read_array.remote(file1)
id2 = read_array.remote(file2)
id = add.remote(id1, id2)
sum = ray.get(id)
```

- Blue variables are Object IDs
- Similar to futures



Return `id1` immediately, before `read_array()` finishes

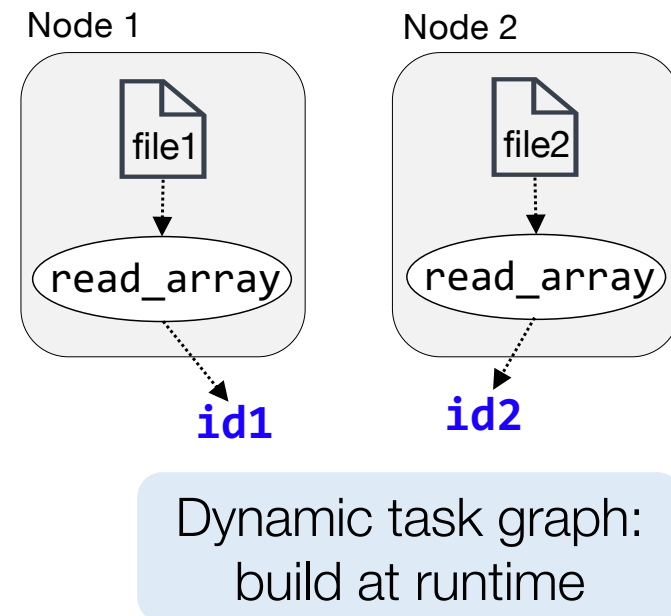
## Task API

```
@ray.remote
def read_array(file):
    # read ndarray "a"
    # from "file"
    return a
```

```
@ray.remote
def add(a, b):
    return np.add(a, b)
```

```
id1 = read_array.remote(file1)
id2 = read_array.remote(file2)
id = add.remote(id1, id2)
sum = ray.get(id)
```

- Blue variables are Object IDs
- Similar to futures





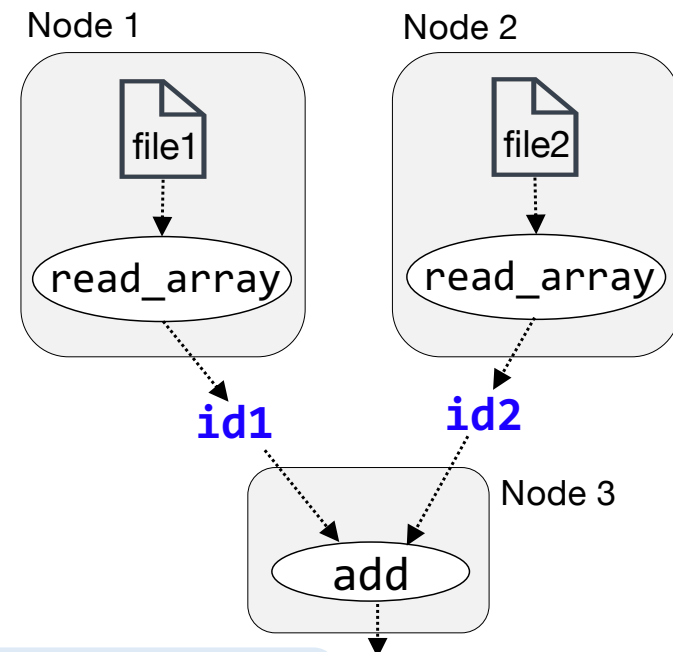
# Task API

```
@ray.remote  
def read_array(file):  
    # read ndarray "a"  
    # from "file"  
    return a
```

```
@ray.remote  
def add(a, b):  
    return np.add(a, b)
```

```
id1 = read_array.remote(file1)  
id2 = read_array.remote(file2)  
id = add.remote(id1, id2)  
sum = ray.get(id)
```

- Blue variables are Object IDs
- Similar to futures



Every task scheduled, but not finished yet

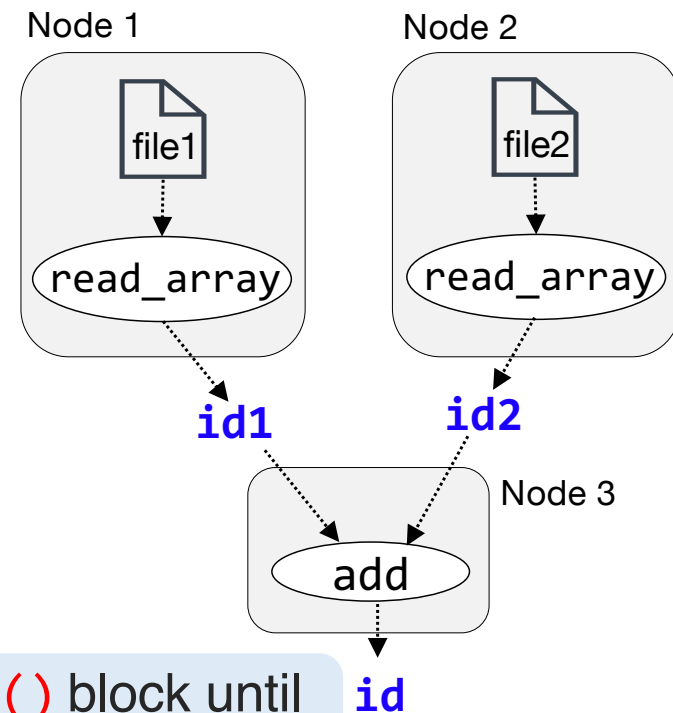
# Task API

```
@ray.remote  
def read_array(file):  
    # read ndarray "a"  
    # from "file"  
    return a
```

```
@ray.remote  
def add(a, b):  
    return np.add(a, b)
```

```
id1 = read_array.remote(file1)  
id2 = read_array.remote(file2)  
id = add.remote(id1, id2)  
sum = ray.get(id)
```

- Blue variables are Object IDs
- Similar to futures



`ray.get()` block until result available

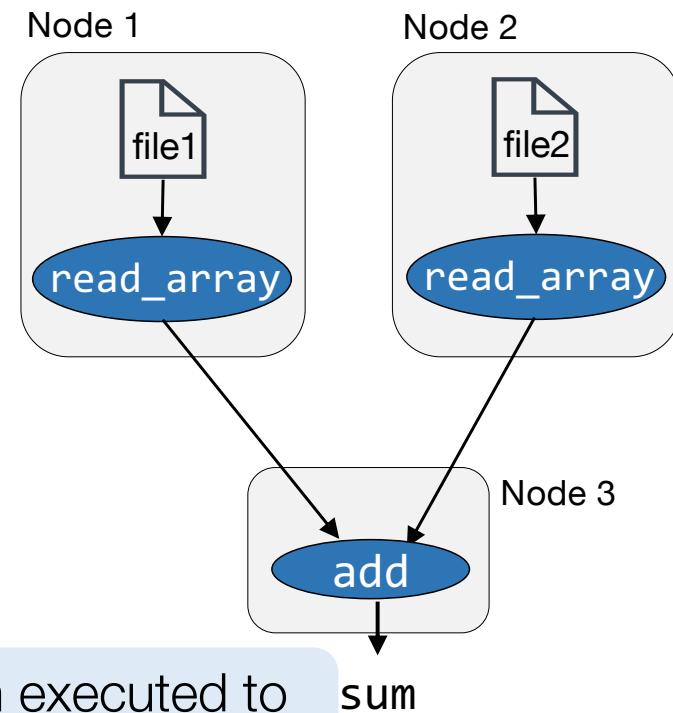
# Task API

```
@ray.remote  
def read_array(file):  
    # read ndarray "a"  
    # from "file"  
    return a
```

```
@ray.remote  
def add(a, b):  
    return np.add(a, b)
```

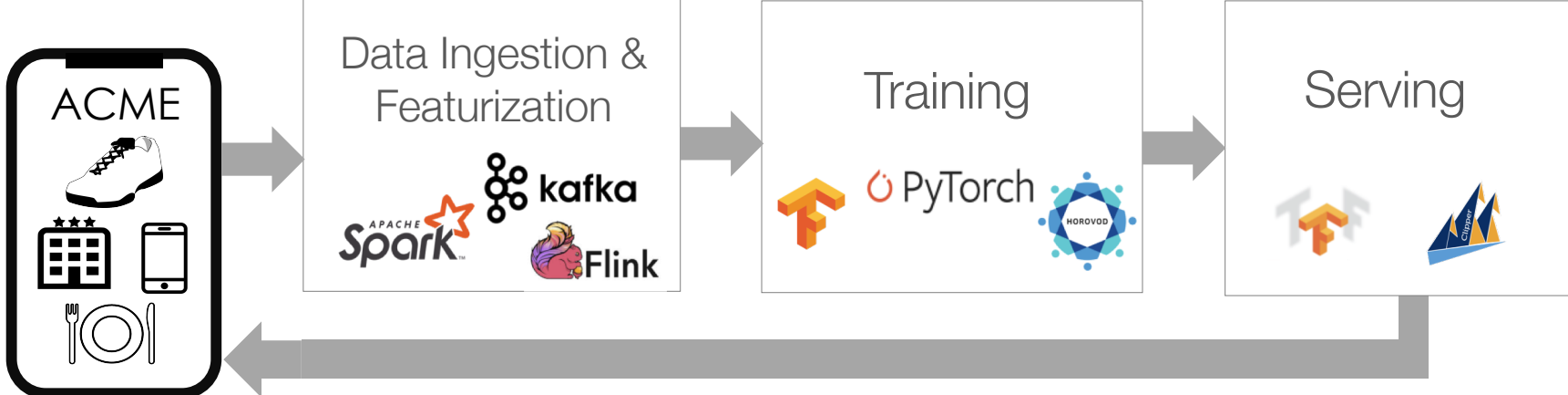
```
id1 = read_array.remote(file1)  
id2 = read_array.remote(file2)  
id = add.remote(id1, id2)  
sum = ray.get(id)
```

- Blue variables are Object IDs
- Similar to futures



Task graph executed to compute sum

# Example: In-app promotion



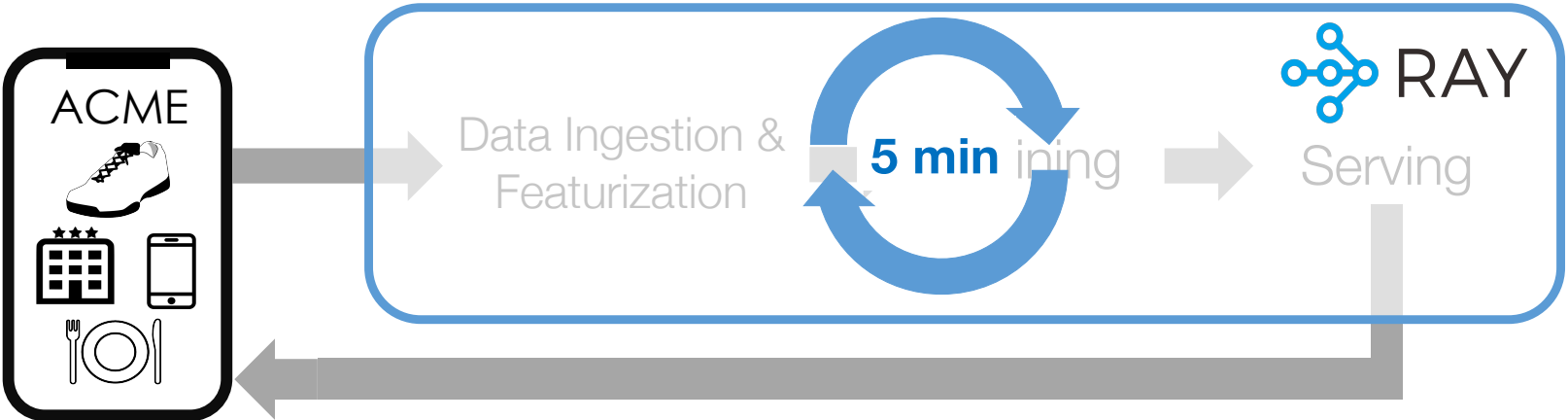
Model updated every **1 day**

↓ + 5% CTR

Model updated every **1 hour** (using state-of-the-art solution)

↓  
?

# Ray: unified platform for distributed apps



Model updated every **1 day**

↓ **+ 5%** CTR

Model updated every **1 hour** (using state-of-the-art solution)

↓ **+ 1%** CTR

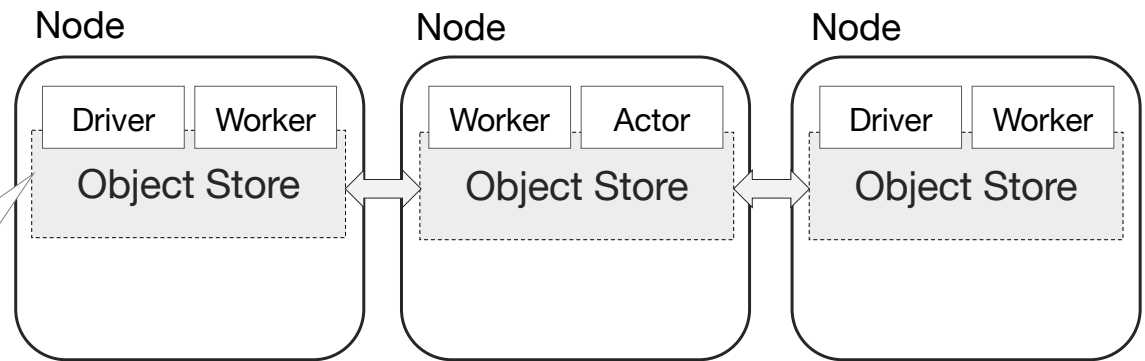
Model updated every **5 min** using **Ray**

# Ray Architecture

## In-memory obj. store

- Immutable objects

Serialization using  
Apache Arrow

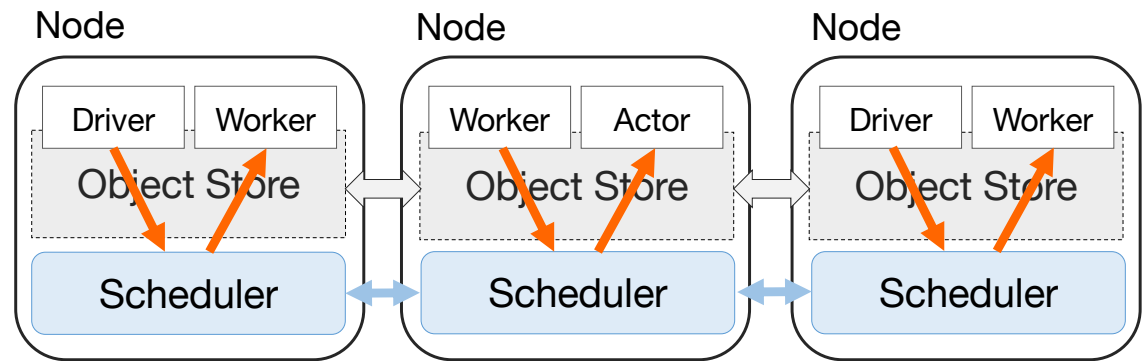


# Ray Architecture

## In-memory obj. store

- Immutable objects

## Distributed scheduler

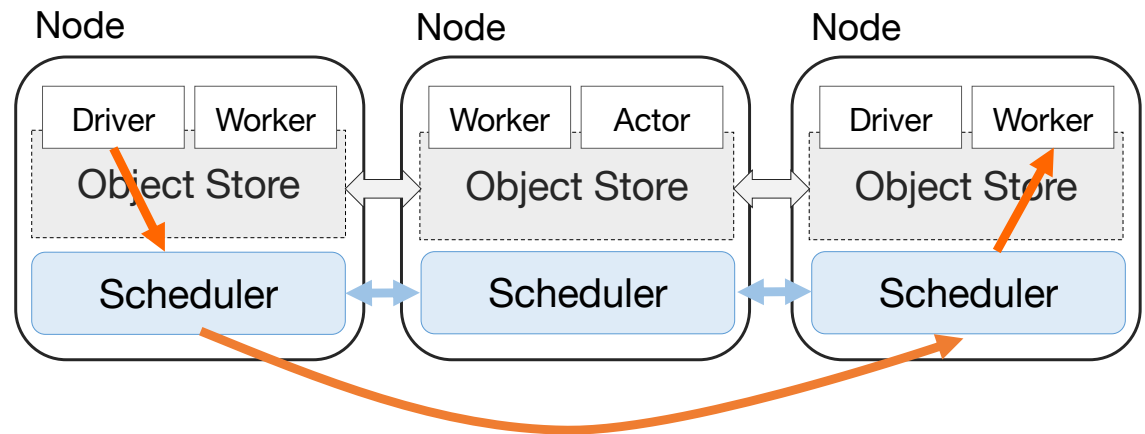


# Ray Architecture

## In-memory obj. store

- Immutable objects

## Distributed scheduler





# Ray Architecture

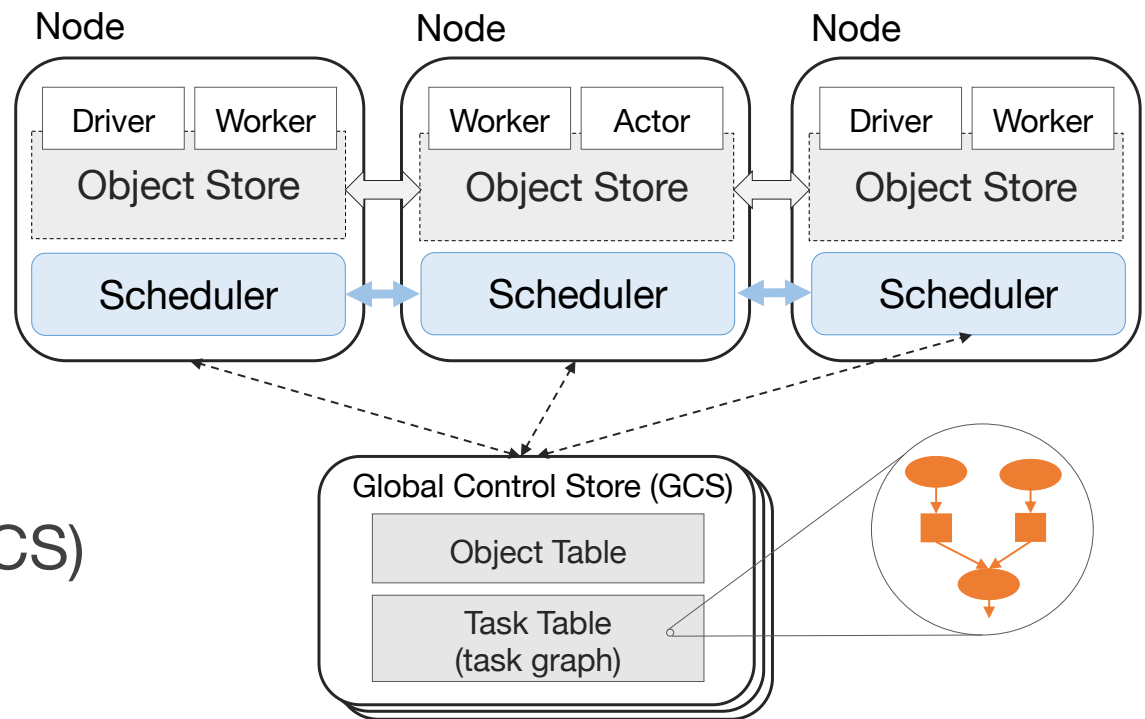
## In-memory obj. store

- Immutable objects

## Distributed scheduler

## Central control store (GCS)

- Stateless components



# Ray Architecture

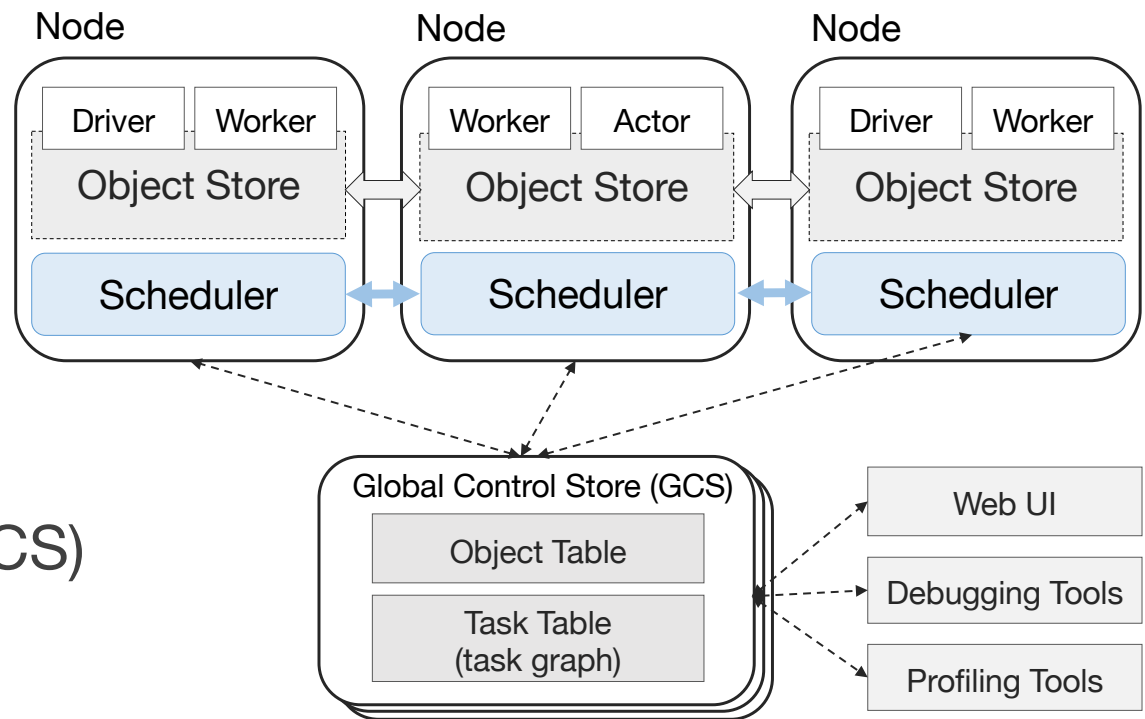
## In-memory obj. store

- Immutable objects

## Distributed scheduler

## Central control store (GCS)

- Stateless components



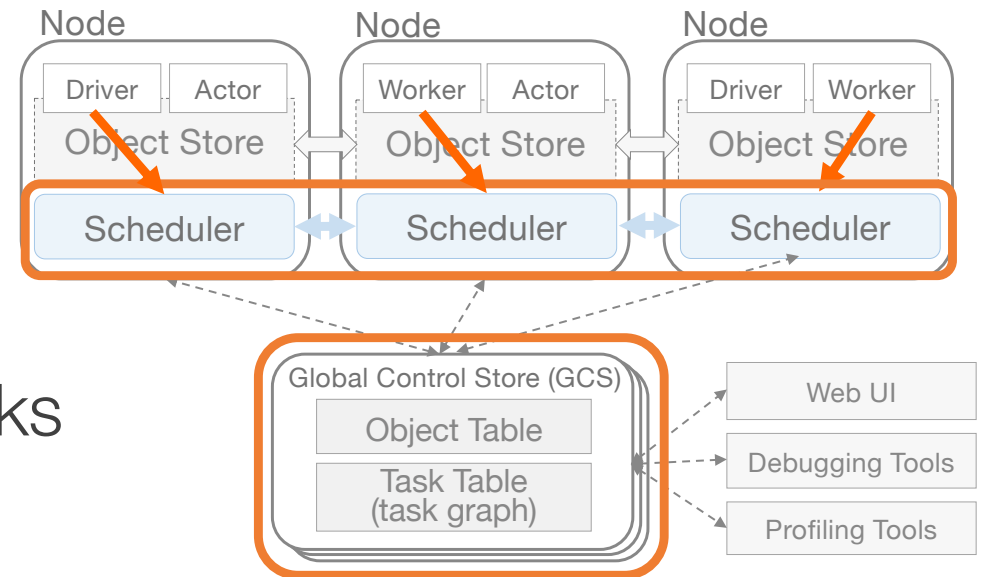
# Scalability

Decentralized scheduler

Sharded GCS

Any worker can submit tasks

- Driver not a bottleneck



# Scalability

Decentralized scheduler

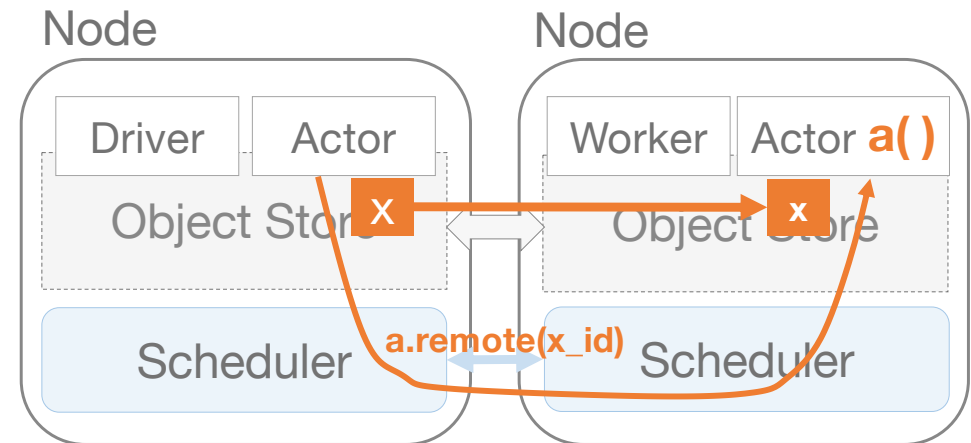
Sharded GCS

Any worker can submit tasks

- Driver not a bottleneck

Actors arguments sent by reference instead inline

- Avoid unnecessary copies
- Can optimize transfers (e.g., parallel transfers, multicast)



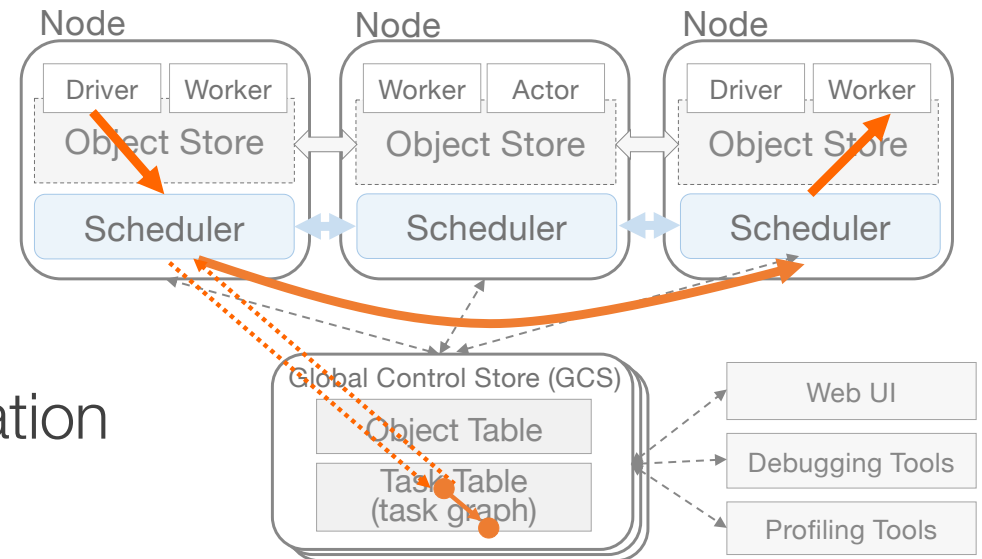
# Fault Tolerance

## Lineage based

- Replay computation to reconstruct lost objects

## Classic solution

- Store lineage before invocation
- Insert delay on critical path



# Fault Tolerance

## Lineage based

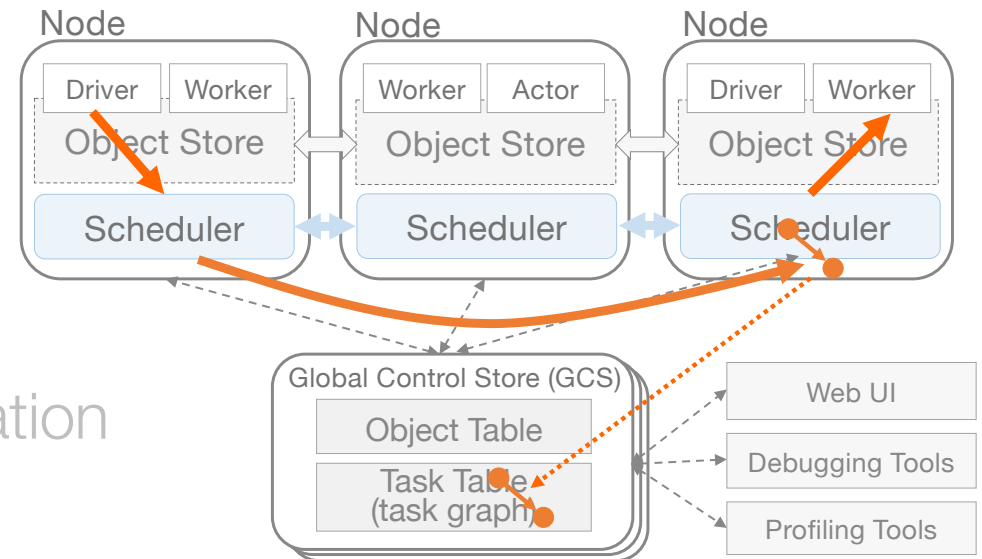
- Replay computation to reconstruct lost objects

## Classic solution

- Store lineage before invocation
- Insert delay on critical path

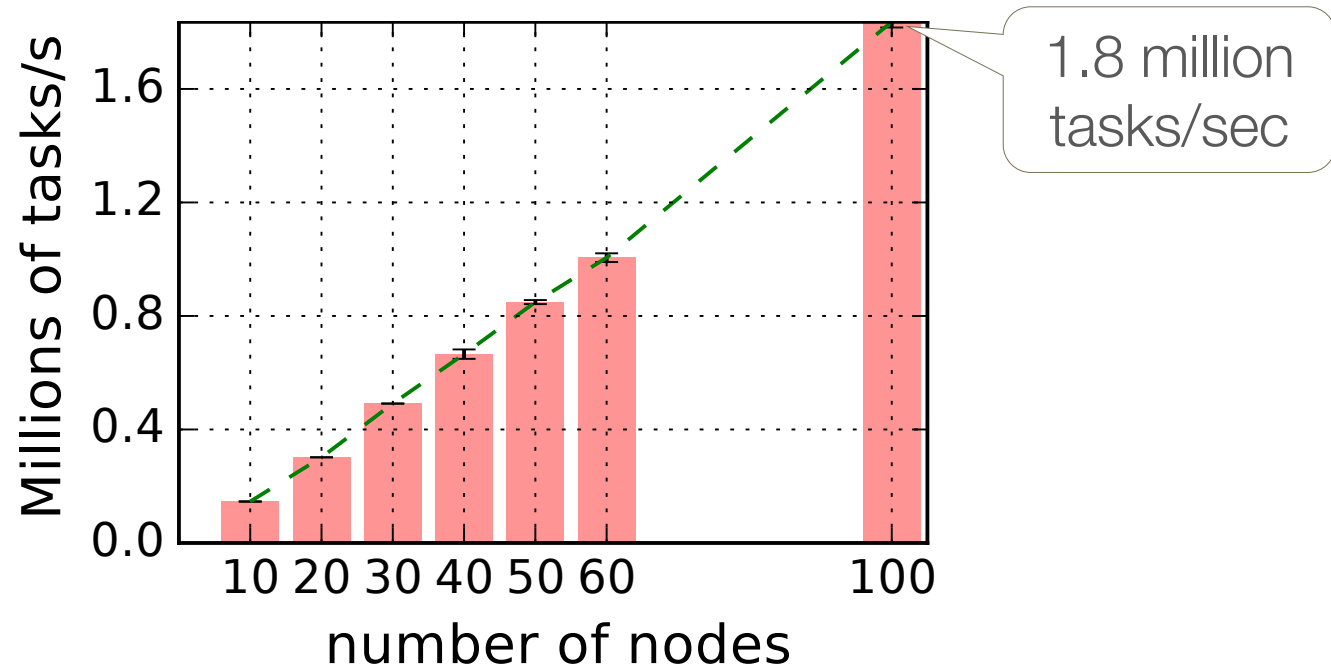
## Lineage stash\*

- Store lineage in task call and flush it



\*"Lineage Stash: Fault Tolerance Off the Critical Path", Stephanie Wang et al, SOSP '19

# Scalability & Performance



Latency of local task execution: ~300 us

Latency of remote task execution: ~1ms

# Ray vs specialized systems



Match performance of specialized systems



## Ongoing work

### Flexible scheduling policies

- Need to support *conflicting* policies, e.g., affinity, antiaffinity, locality, gang scheduling, ...

### Improved garbage collection for object store

- Currently LRU, but not "good enough"
- Ideally, global reference counting

RL Library



Hyperparam.  
Search



Data  
Processing



Distributed  
applications



RAY

General-purpose distributed computing  
framework for Python (and Java)

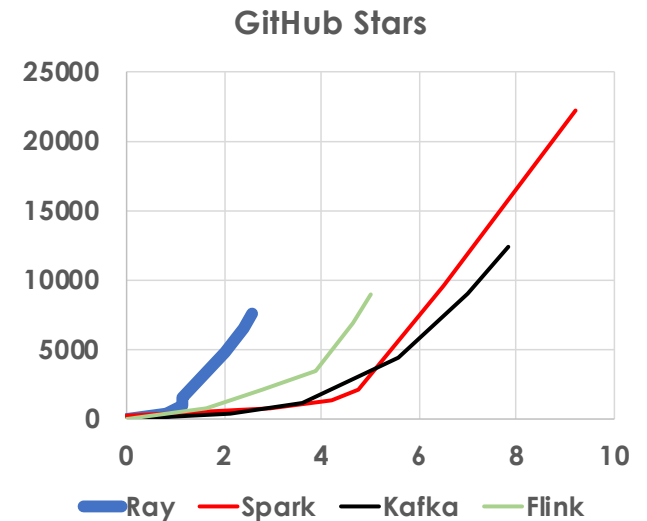
# Growing adoption



200 contributors from  
40+ companies

Sold out tutorials at O'Reilly AI

Included in AWS Sage Maker



J.P.Morgan

Morgan Stanley

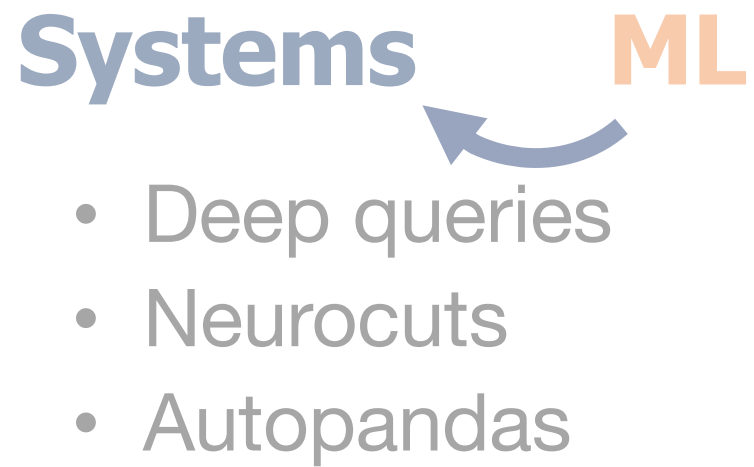


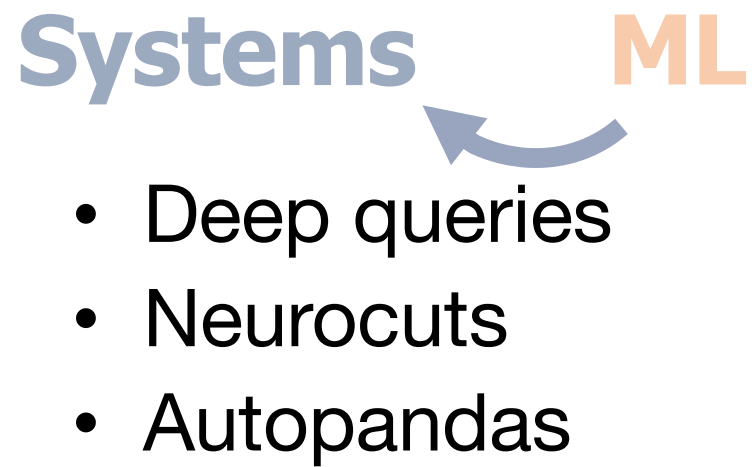
PRIMER



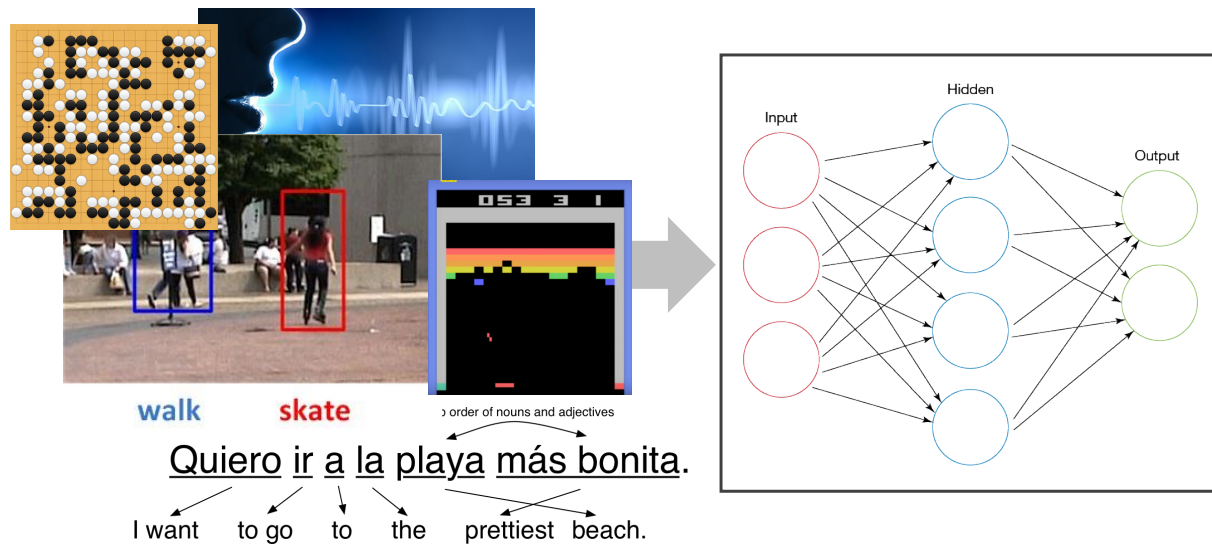
Microsoft







# “Classic” DL/RL apps

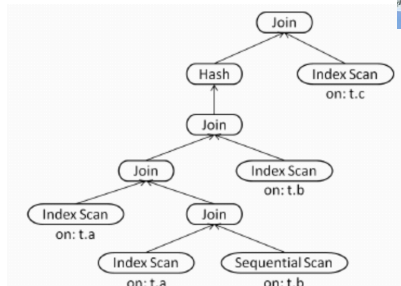
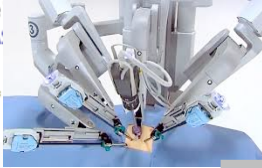


- Speech recognition
- Video recognition
- Language translation
- ...

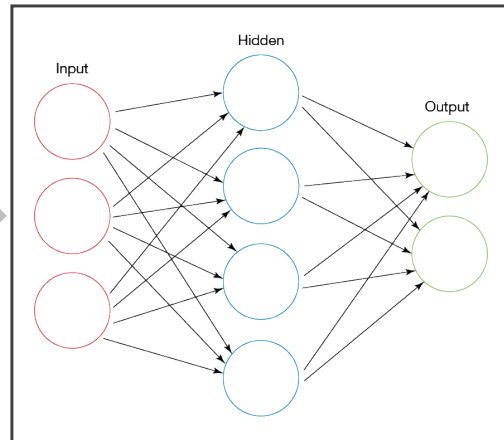
Human tasks: 100% accuracy not expected

# Systems problems

```
SELECT COALESCE(users.state, '') AS "_g1",  
       users.state AS `users.state`  
       COUNT(DISTINCT orders.id) AS  
FROM orders  
LEFT JOIN users ON orders.user_id =
```



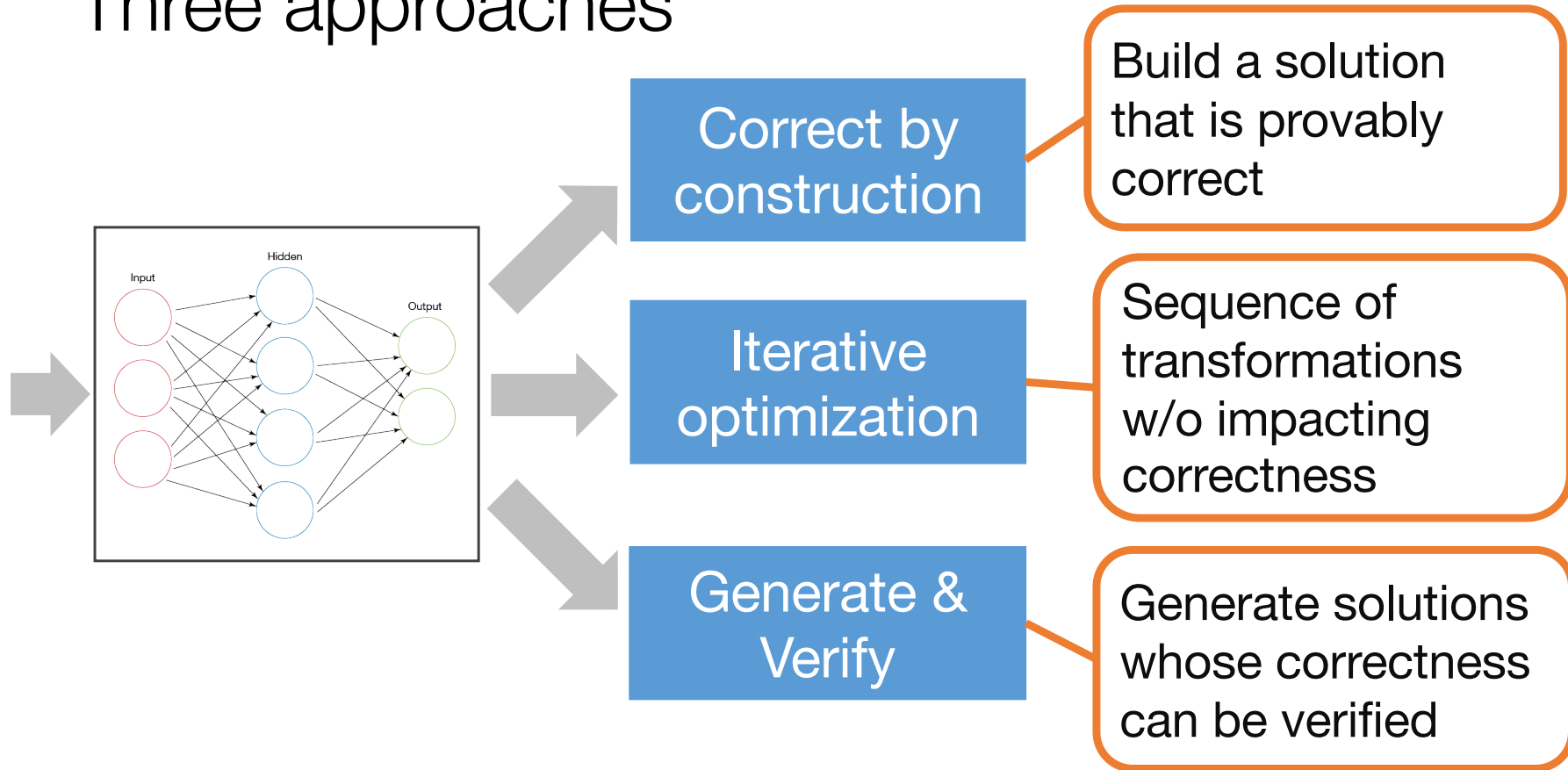
**REJECTED**  
**MORTGAGE**  
**APPLICATION**



Program synthesis  
Mortgage decisions  
Robotic surgery  
...

Must ensure correctness and explainability!

# Three approaches





# Database joins

Calculate tax owned by  
“Manager I” employees

```
SELECT SUM(sal.salary*tax.rate)
FROM emp, sal, tax
WHERE emp.position = sal.position AND
      tax.country = sal.country AND
      emp.position = 'Manager I'
```



emp_id	position	country
1	Manager II	USA
2	Engineer I	CAN
3	Engineer II	USA
...	...	..

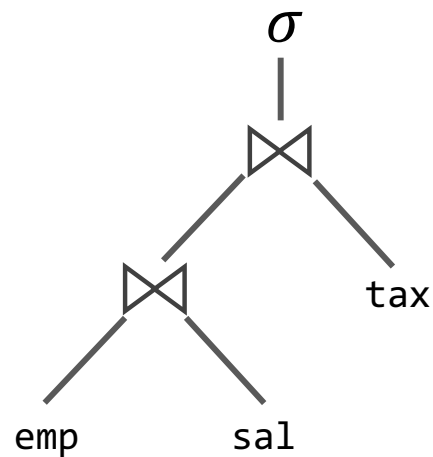
sal_id	position	salary
1	Manager I	120000.00
2	Manager II	150000.00
3	Engineer I	78000.00
4	Engineer II	91000.00

tax_id	country	rate
1	USA	0.32
2	CAN	0.45
3	CHN	0.17
4	...	...

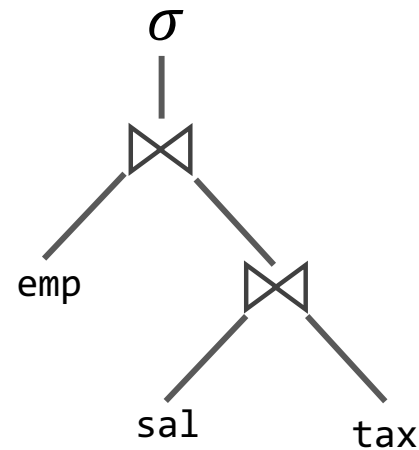
# Join Optimization

```
SELECT SUM(sal.salary*tax.rate)
FROM emp, sal, tax
WHERE emp.position = sal.position AND
      tax.country = sal.country AND
      emp.position = 'Manager I'
```

In what order do you perform joins?



cost = 1200

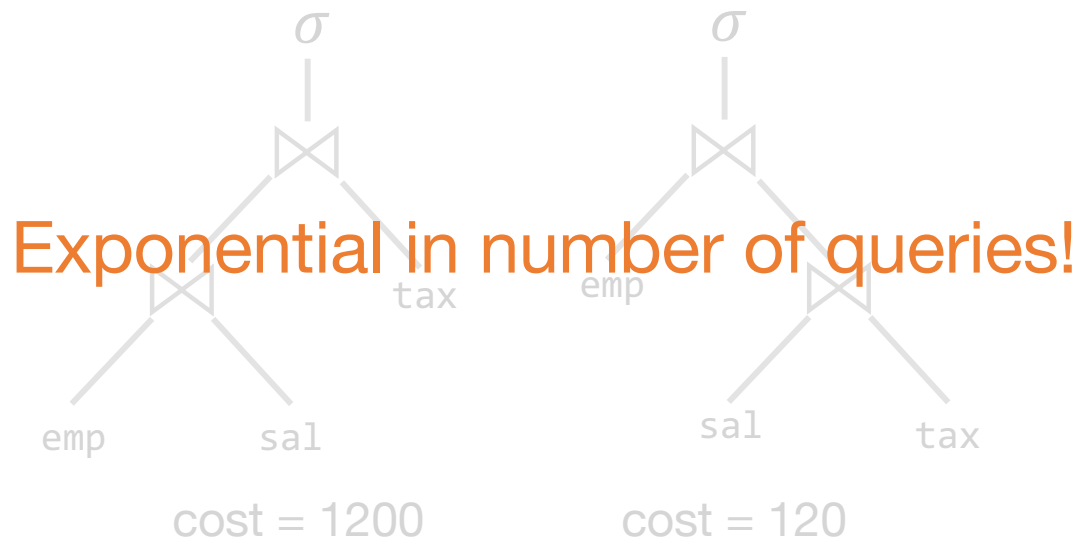


cost = 120

# Join Optimization

```
SELECT SUM(sal.salary*tax.rate)
FROM emp, sal, tax
WHERE emp.position = sal.position AND
      tax.country = sal.country AND
      emp.position = 'Manager I'
```

In what order do you perform joins?



# 40 years of heuristics!

**Left-deep:** maximize index usage

**Right-deep:** maximize hash table re-use

**Zig-zag:** union of LD and RD

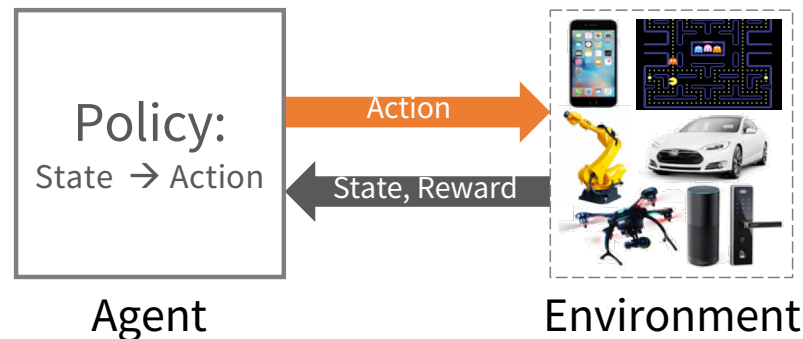
**Greedy:** exploit linear cost models

**IK-KBZ:** exploit Star Schemas

**GEQO:** genetic algorithms in Postgres

Can ML replace programmed heuristics  
with efficient and data-driven strategies?

# Reinforcement Learning (RL)



Agent continually learning by interacting with env.  
Compute policy (i.e., state  $\rightarrow$  action) to maximize reward

# Deep Query\*: Join plans with Deep RL



## Why reinforcement learning (RL)?

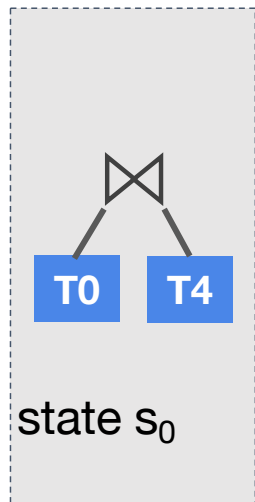
- Natural formulation
  - **State**: set of tables joined so far
  - **Action**: table to join next
  - **Reward**: negative of estimated cost
- Use Deep Q-Learning

\*"Learning to Optimize Join Queries With Deep Reinforcement Learning", Sanjay Krishnan et al  
(<https://arxiv.org/abs/1808.03196>)

# Deep Query: Build join plans with Deep RL

Why reinforcement learning (RL)?

- Natural formulation

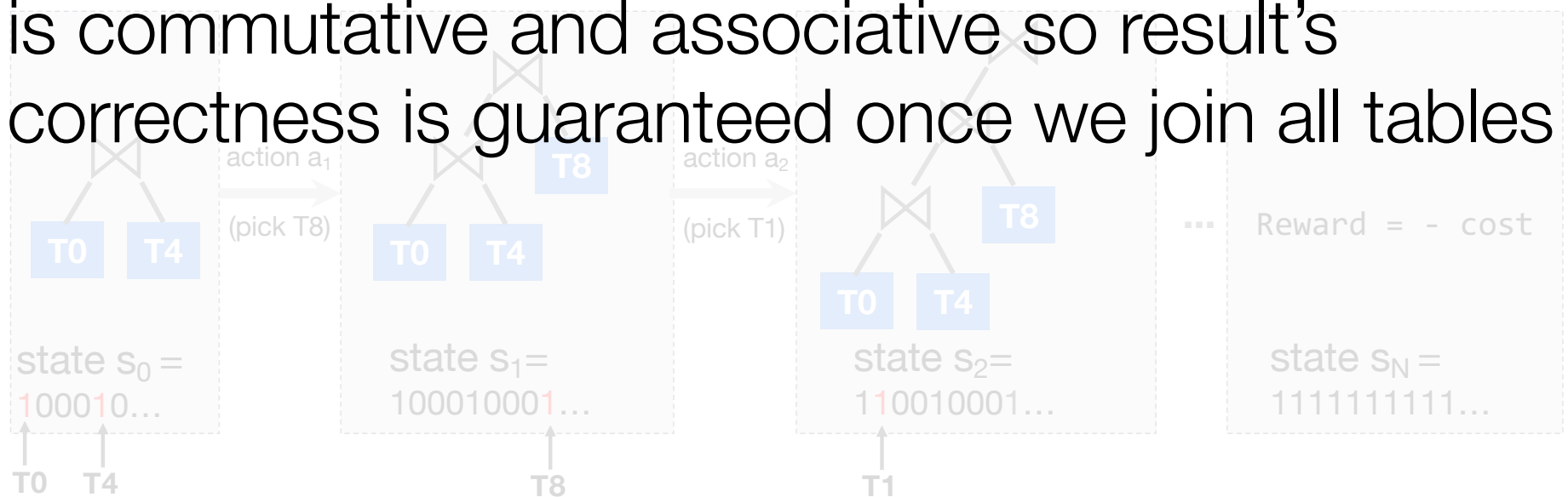


# Deep Query: Build join plans with Deep RL

Why reinforcement learning (RL)?

Natural formulation

**Correct by construction:** the "join" operation is commutative and associative so result's correctness is guaranteed once we join all tables





## Deep Query: Properties

**Generalize** to unseen queries

**Adapt** to workload and hardware characteristics

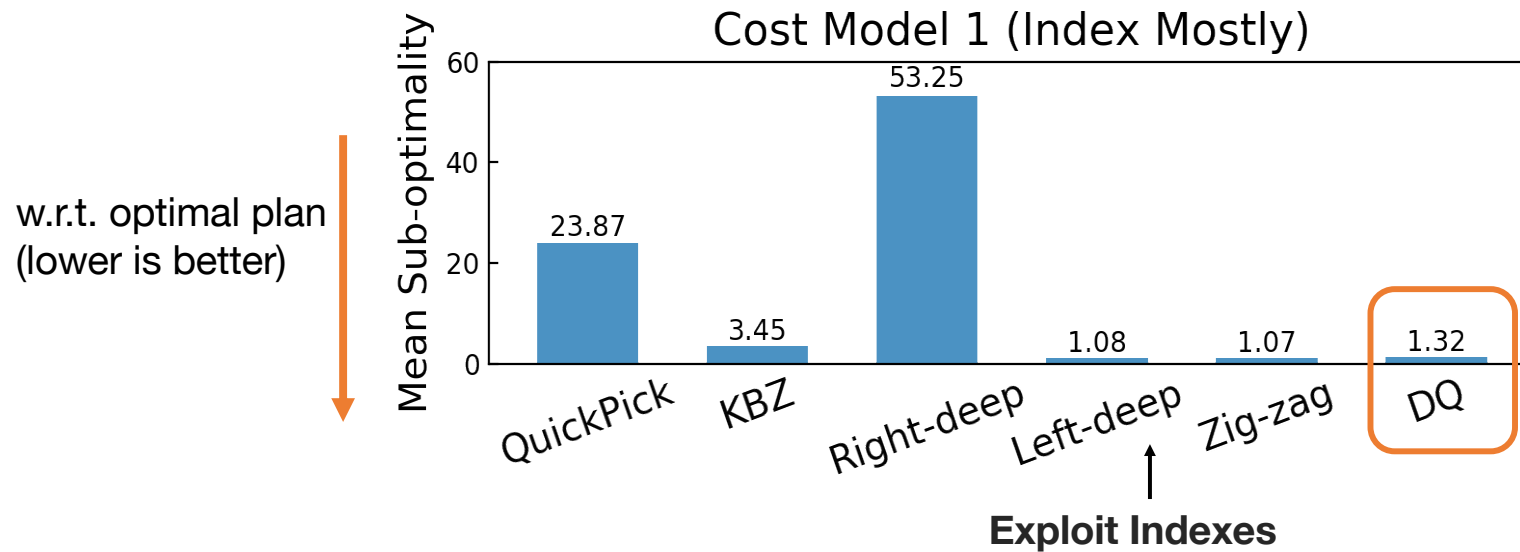
**Efficient** planning (order of magnitude faster)

113 queries, IMDb dataset, 21 tables

4–17 way joins (avg ~8 relations per query)

Cost models:

- Model 1: Index Mostly
- Model 2: Hybrid Hash



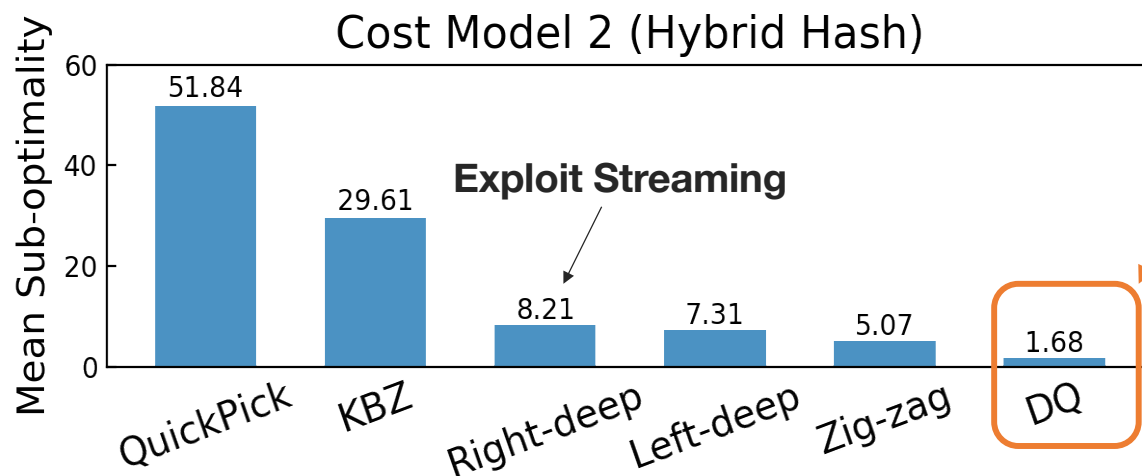
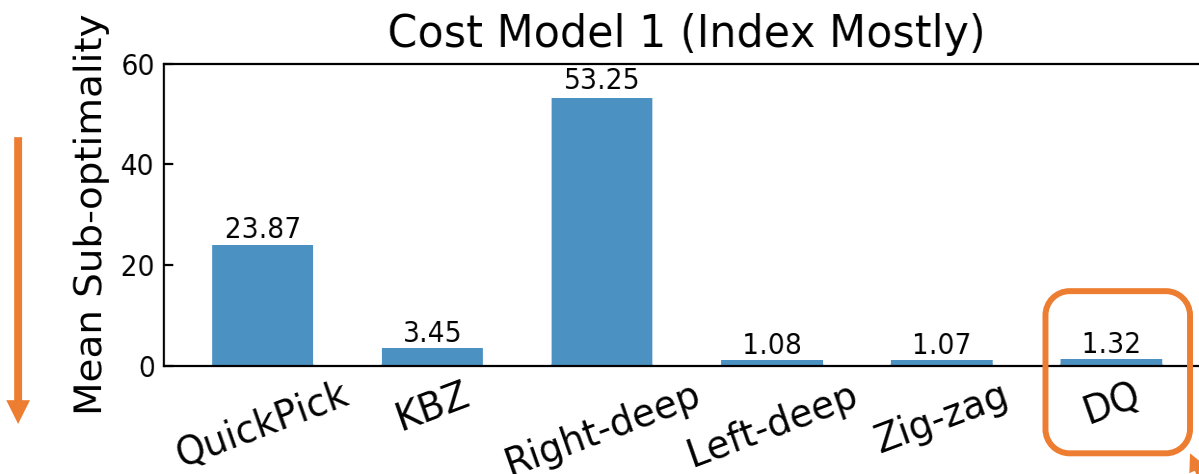
113 queries, IMDb dataset, 21 tables  
 4–17 way joins (avg ~8 relations per query)

Cost models:

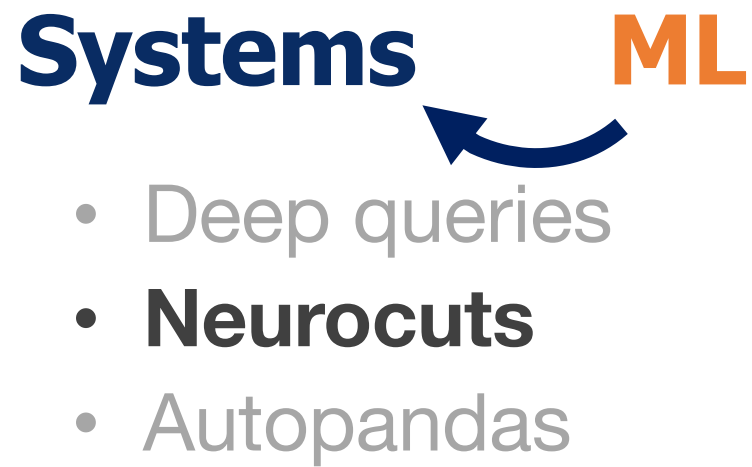
- Model 1: Index Mostly
- Model 2: Hybrid Hash



w.r.t. optimal plan  
(lower is better)



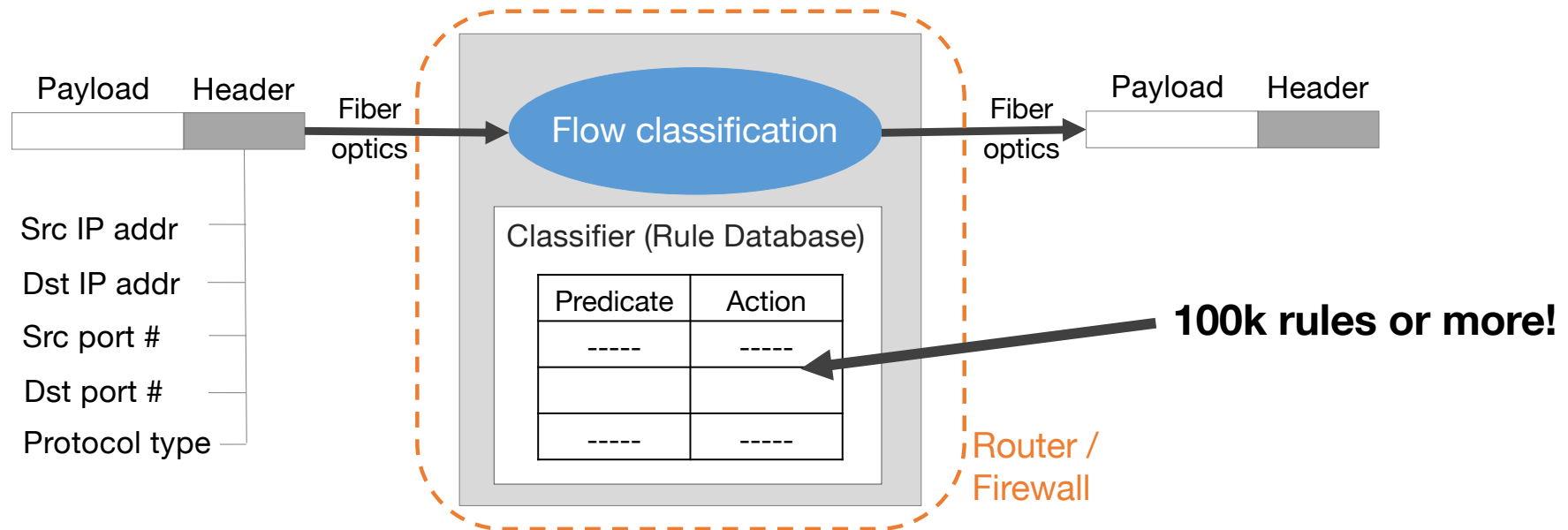
Robust



# Packet Classification

## Fundamental problem in networking

- Building block for access control, QoS, defense against attacks

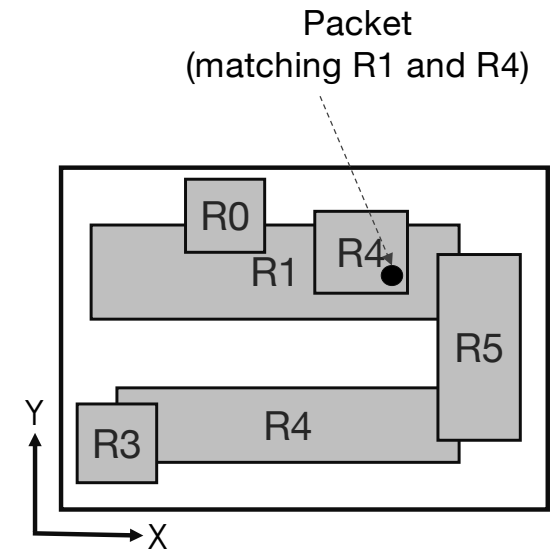


# The problem

Similar to point-location in a hypercube

Hard time-space tradeoff:

- $O(\log N)$  time and  $O(N^d)$  space
- $O(\log^{d-1} N)$  time and  $O(N)$  space
- $N$ : # of rules;  $d$ : # of attributes
  - In our case:  $N \approx 100K$ ,  $d = 5$



But harder: rules overlap and have priorities

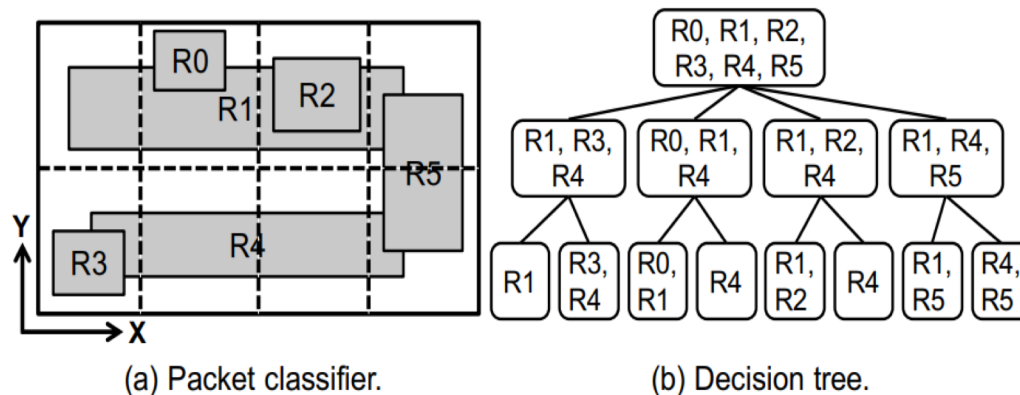
# 20+ years of work

## Hardware

- Expensive and power hungry → prohibitive for large classifiers

## Software

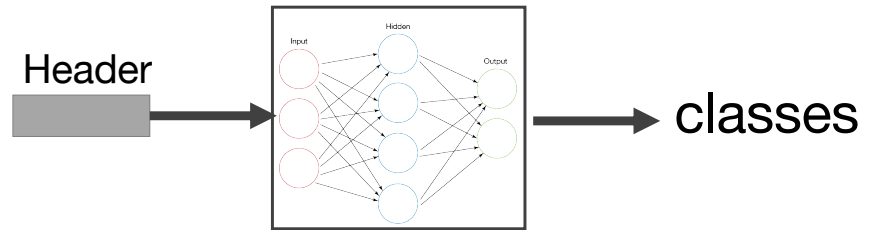
- Build a multi-dimensional "decision tree" – really a k-d index
- HiCuts ('99), HyperCuts ('03), EffiCuts ('10), CutSplit ('15)
  - All rely on hand-tuned heuristics, which are brittle and not optimal





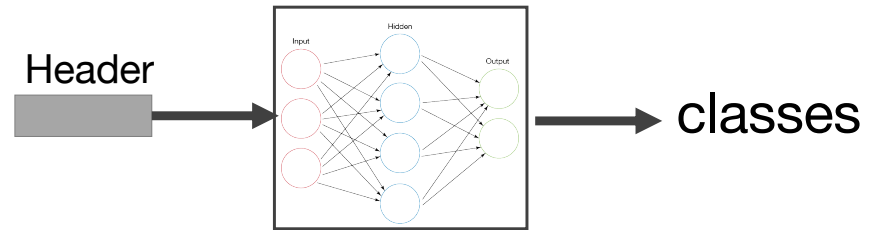
# Two approaches

## 1. End-to-end solution

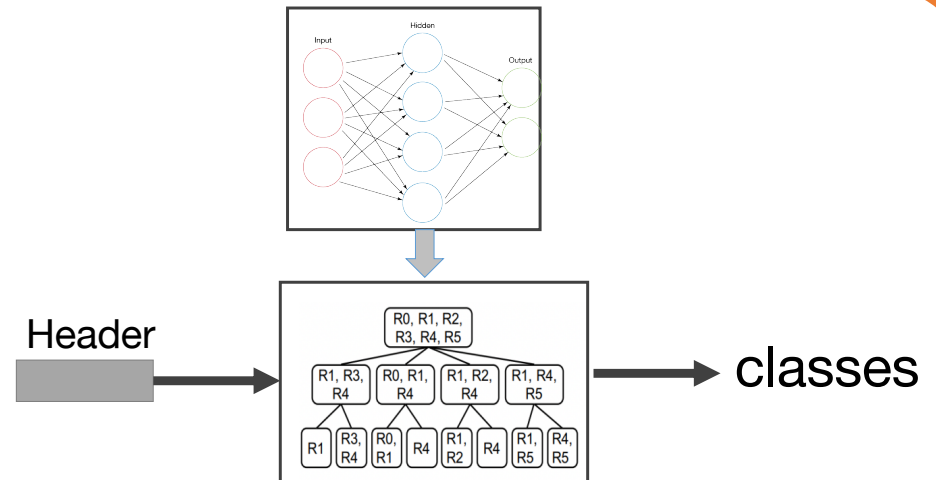


# Two approaches

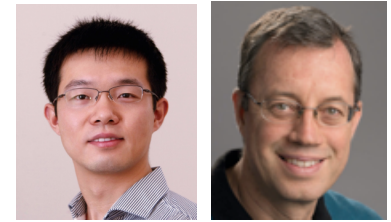
1. End-to-end solution



2. Build classification tree



# NeuroCuts\*: Building decision trees with Deep RL



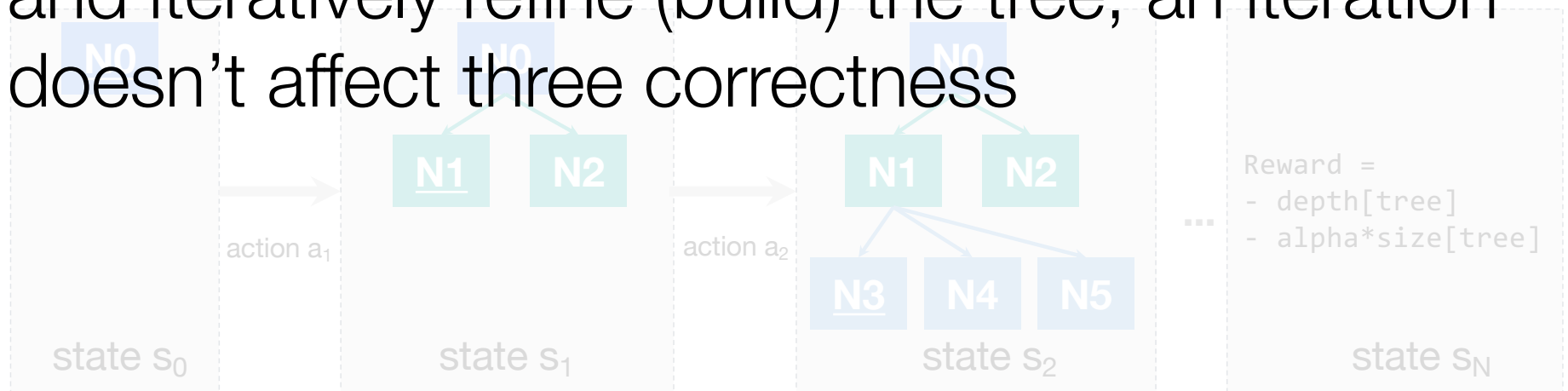
\*"Neural Packet Classification", Eric Liang et al, SIGCOMM 2019

# NeuroCuts: Building decision trees with Deep RL

Why reinforcement learning (RL)?

- Simulation is reality: build a tree in simulation → deploy to

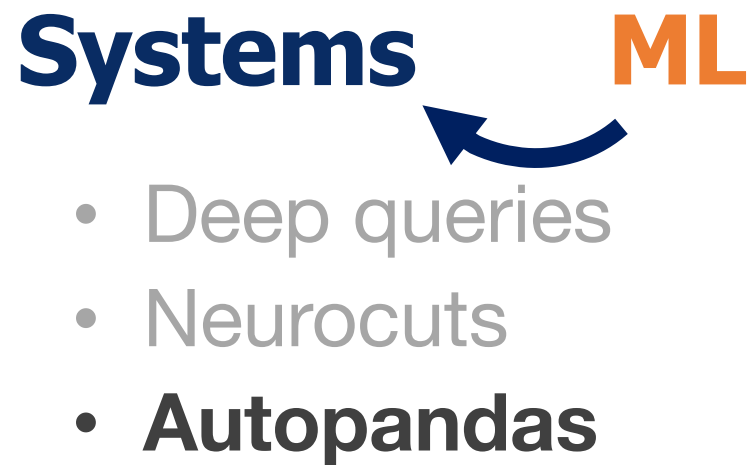
**Iterative optimization:** start from a single-node and iteratively refine (build) the tree; an iteration doesn't affect three correctness



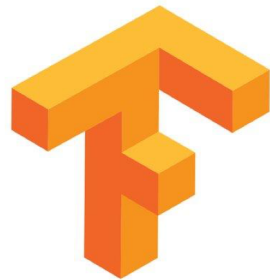
# Results

Classification time: median **18%** faster than state-of-the-art

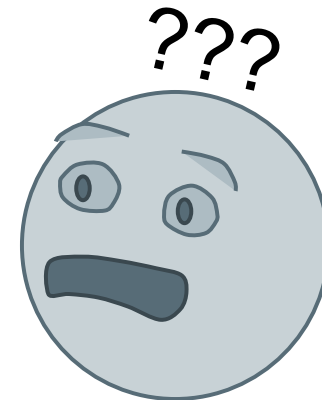
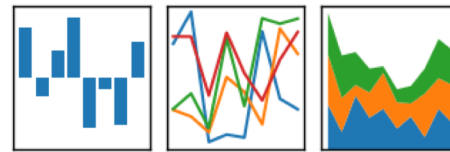
**3x** better either in space or time than any previous solution



# API Explosion!



pandas  
 $y_{it} = \beta' x_{it} + \mu_i + \epsilon_{it}$



# How to cope? *StackOverflow*



How do I turn this:

	weight	
	kg	lbs
cat	1	2
dog	2	4

into this:

		weight
cat	kg	1
	lbs	2
dog	kg	2
	lbs	4

in pandas?



# Problems with *StackOverflow*



How do I turn this:

	weight	
	kg	lbs
cat	1	2
dog	2	4

into this:

		weight
cat	kg	1
	lbs	2
dog	kg	2
	lbs	4

in pandas?

**Inefficient  
Solutions**

Well, you need to  
start by building the  
index  
`pd.MultiIndex(...`



# Problems with *StackOverflow*



How do I turn this:

	weight	
	kg	lbs
cat	1	2
dog	2	4

into this:

		weight
cat	kg	1
	lbs	2
dog	kg	2
	lbs	4

in pandas?

**Inefficient  
Solutions**

**4 days later!**

Just use the  
stack function



# Goal: *StackOverflow* for APIs via Program Synthesis



How do I turn this:

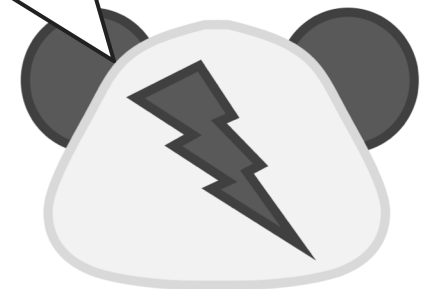
	weight	
	kg	lbs
cat	1	2
dog	2	4

into this:

		weight
cat	kg	1
	lbs	2
dog	kg	2
	lbs	4

in pandas?

```
output = input.stack(  
    level=[1],  
    dropna=True  
)
```

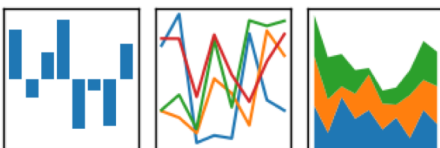


AutoPandas

# Target API: *pandas* (DataFrame transformations)

# pandas

$$y_{it} = \beta' x_{it} + \mu_i + \epsilon_{it}$$

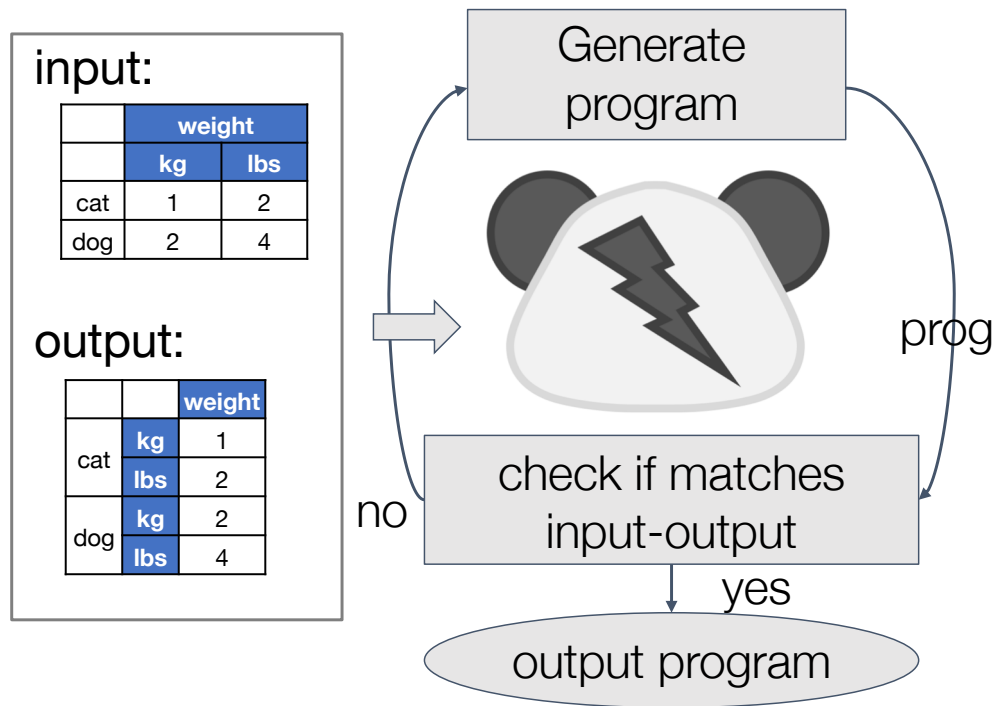


## Premier library for data scientists

A collection of cheat sheets for pandas DataFrame operations, categorized into: General functions, Data manipulations, Missing data handling, Computation, and Function application, Groupby & Window. The cheat sheets list various methods like .isin(), .where(), .mask(), .query(), .stack(), .unstack(), .pivot(), .melt(), .groupby(), etc., with brief descriptions of their functionality.

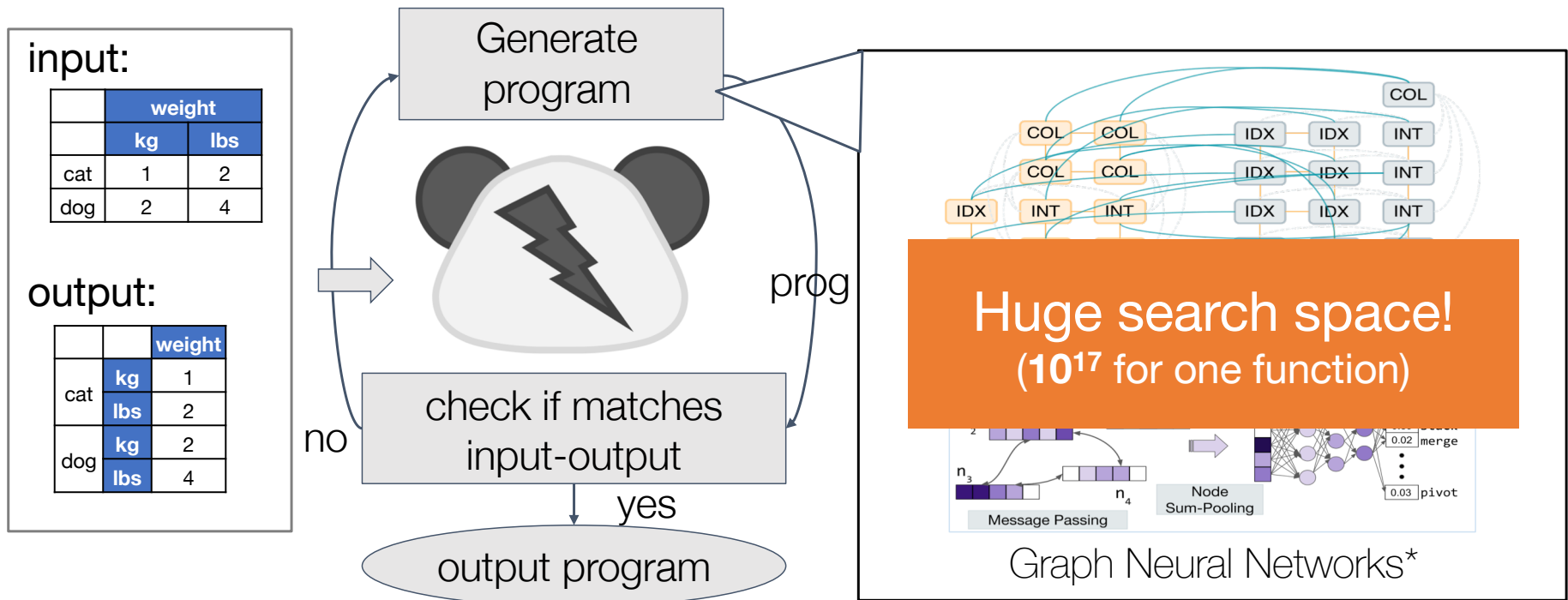
Avg. number of single-function programs: **10<sup>17</sup>** 🤯

# Autopandas\*



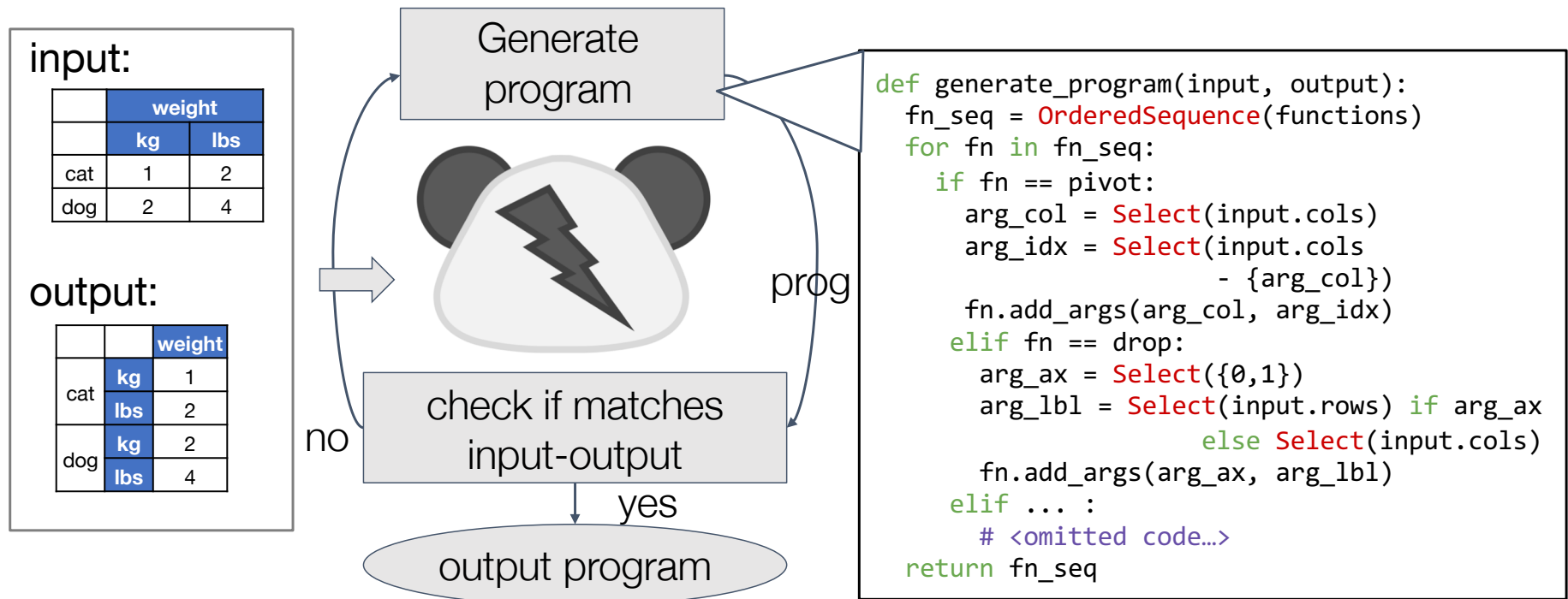
“AutoPandas: Neural-Backed Generators for Program Synthesis”, Rohan Bavishi et al, OOPSLA 2019

# Predict program

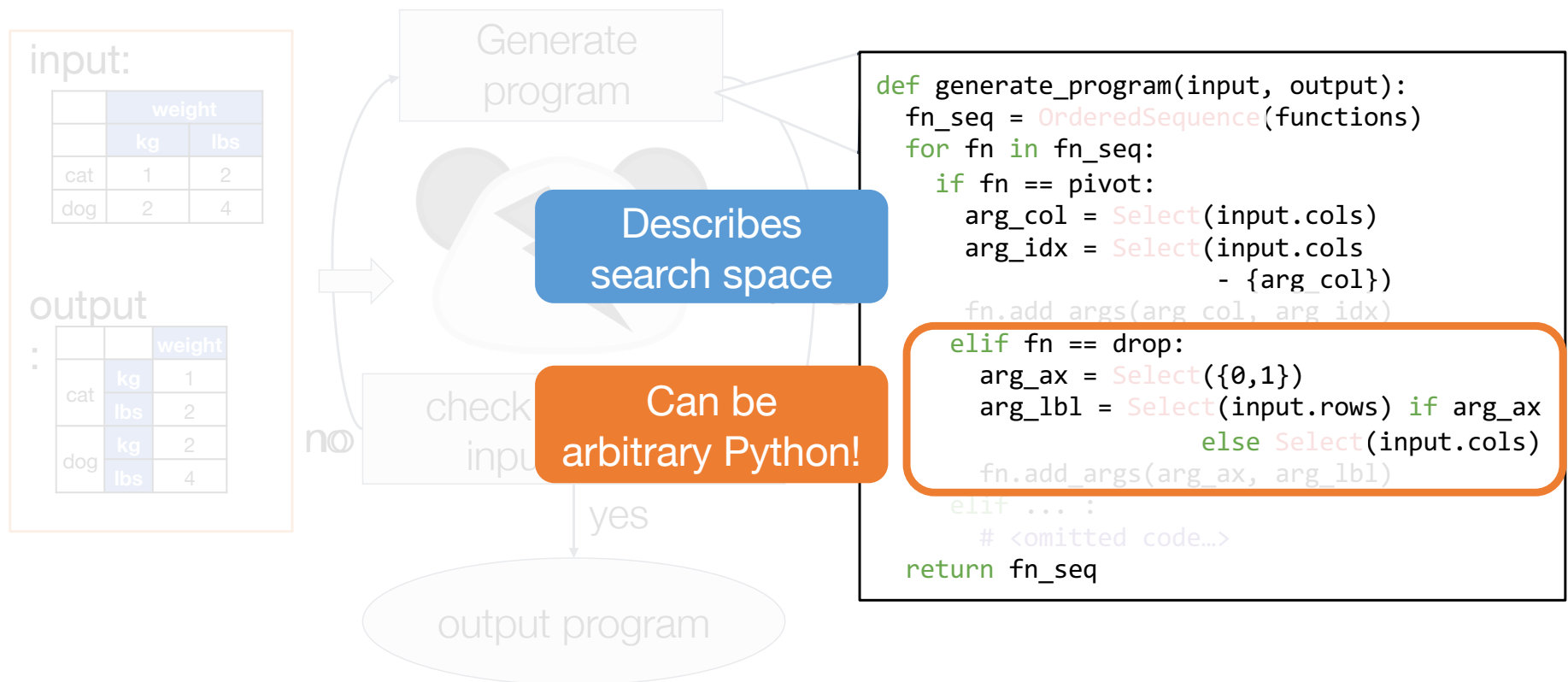


\*M. Allamanis, M. Brockschmidt, and M. Khademi, Learning to represent programs with graphs, ICLR 2018

# Neural-backed program generators

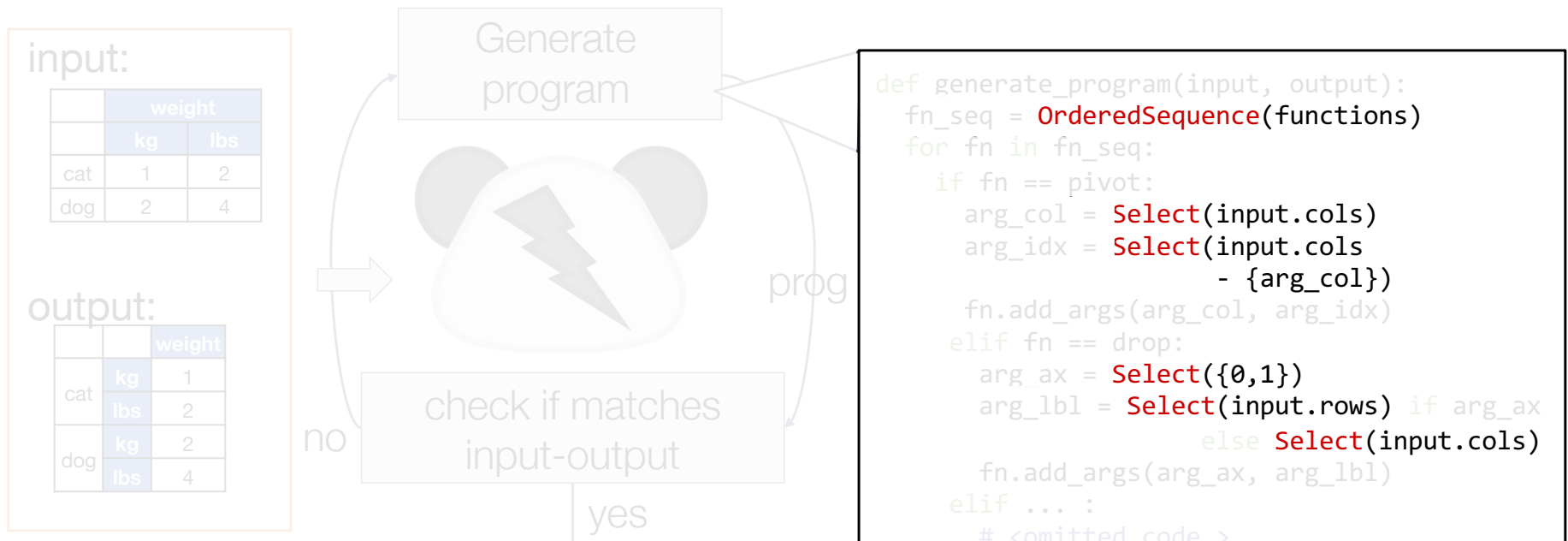


# Neural-backed program generators





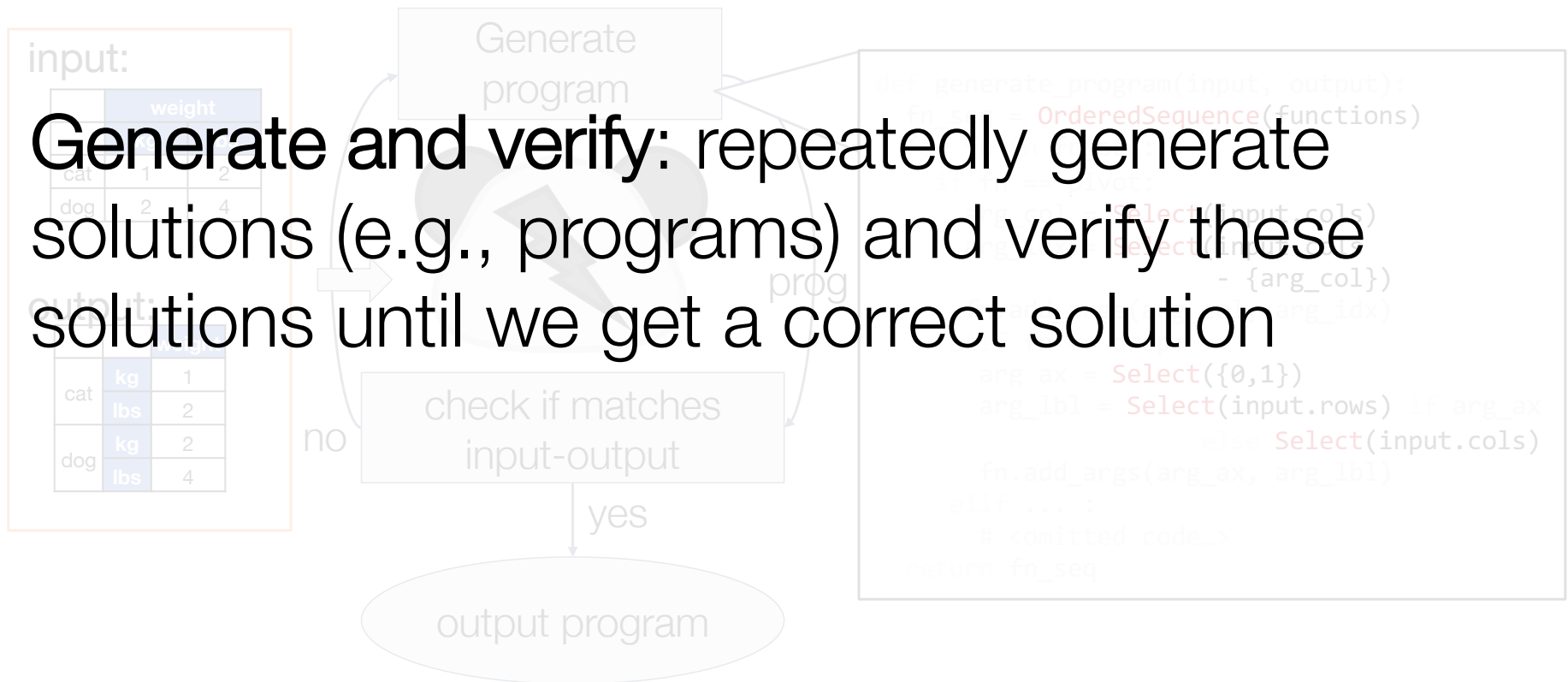
# Neural-backed program generators

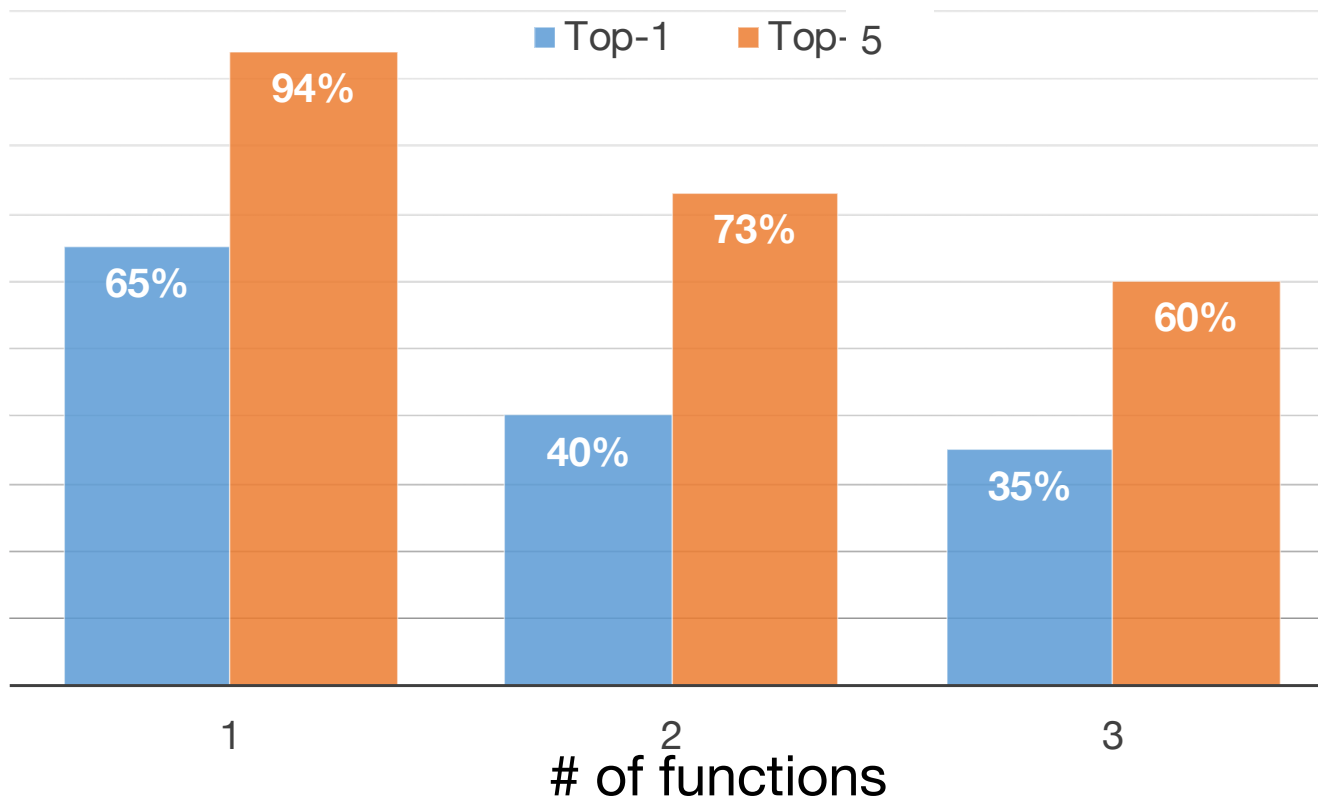


For one function reduce search space by  $10^{12}$   
(from  $10^{17}$  to  $10^5$ )

# Neural-backed program generators

**Generate and verify:** repeatedly generate solutions (e.g., programs) and verify these solutions until we get a correct solution





- 70 benchmarks from StackOverflow, “Python for Data Analysis”, “Data School” videos
- ~95% of code snippets in StackOverflow have length  $\leq 3$  functions

## Summary



**ML** needs **systems**: scale algorithms, easy to program

**Systems** need **ML**: improve state-of-the-art heuristics

Putting the two together → explosive growth

Systems → better ML → better Systems → ...