

Building an Operating System for AI

How Microservices and Serverless Computing Enable
the Next Generation of Machine Intelligence

ALGORITHMIA

Diego Oppenheimer, CEO

diego@algorithmia.com

About Me



Diego Oppenheimer - Founder and CEO - Algorithmia

- Product developer, entrepreneur, extensive background in all things data.
- Microsoft: PowerPivot, PowerBI, Excel and SQL Server.
- Founder of algorithmic trading startup
- BS/MS Carnegie Mellon University

ALGORITHMIA

Make state-of-the-art algorithms
discoverable and **accessible**
to everyone.

Algorithmia.com

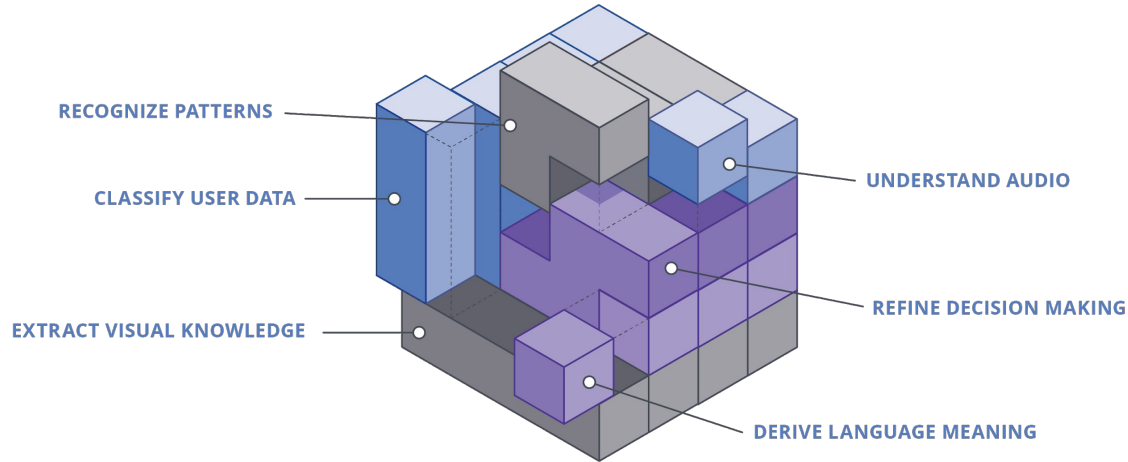
AI/ML scalable infrastructure on demand + marketplace

- Function-as-a-service for Machine & Deep Learning
- Discoverable, live inventory of AI
- Monetizable
- Composable
- **Every developer on earth can make their app intelligent**



“There’s an algorithm for that!”

63K DEVELOPERS 4.8K ALGORITHMS



How do we do it?

- ~4,800 algorithms (50k w/ different versions)
- Each algorithm: 1 to 1,000 calls a second, fluctuates, no devops
- ~15ms overhead latency
- Any runtime, any architecture

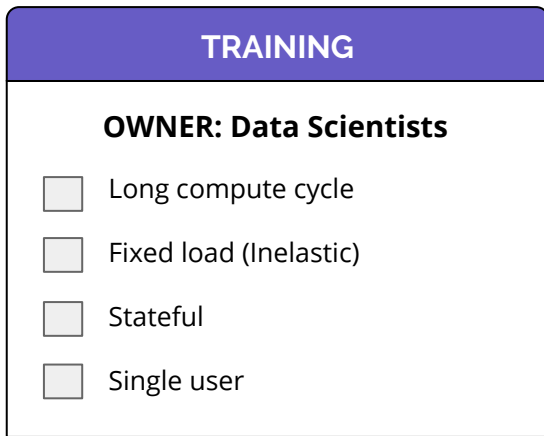
Characteristics of AI

- Two distinct phases: training and inference
- Lots of processing power
- Heterogenous hardware (CPUs, GPUs, TPUs, etc.)
- Limited by compute rather than bandwidth
- “*Tensorflow is open source, scaling it is not.*” - Kenny Daniel

TRAINING

OWNER: Data Scientists

- ☐ Long compute cycle
- ☐ Fixed load (Inelastic)
- ☐ Stateful
- ☐ Single user



Analogous to dev tool chain.

Building and iterating over a model is similar to building an app.

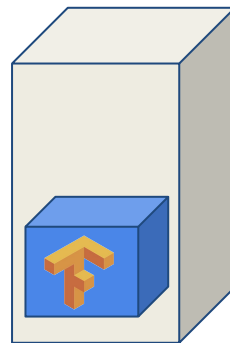
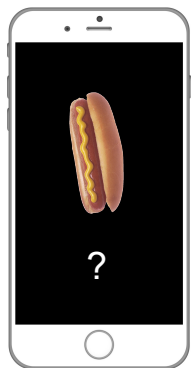
Use Case

Jian Yang made an app to recognize food "SeeFood". Fully trained. Works on his machine.



Use Case

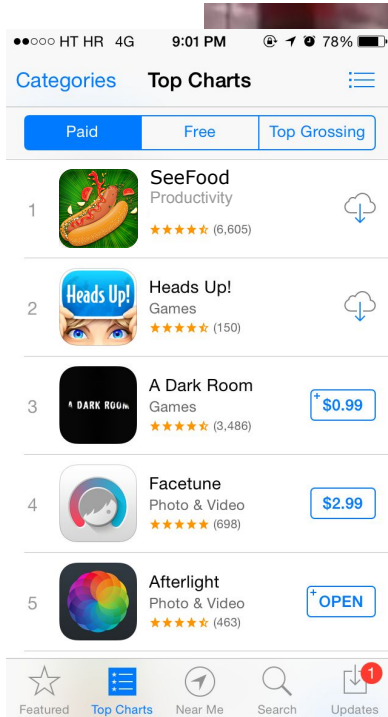
He deployed his trained model to a GPU-enabled server



GPU-enabled
Server

Use Case

The app is a hit!



Use Case

... and now his server is overloaded.



We'll be talking about Microservices & Serverless Computing

MICROSERVICES: the design of a system as independently deployable, loosely coupled services.

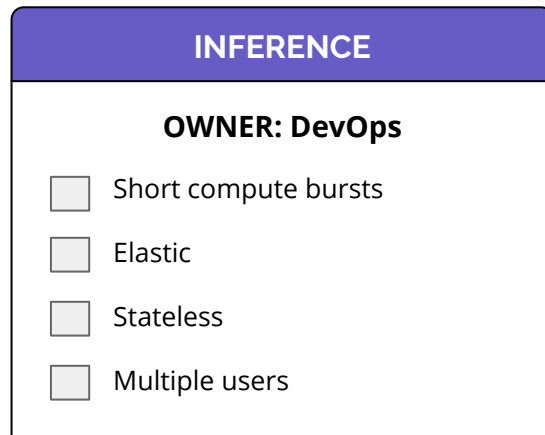
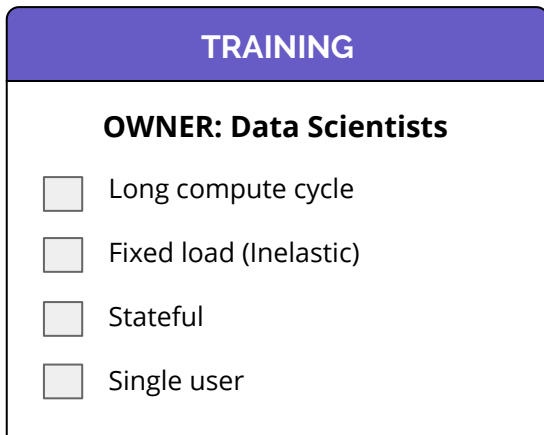
ADVANTAGES

- Maintainability
- Scalability
- Rolling deployments

SERVERLESS: the encapsulation, starting, and stopping of singular functions per request, with a just-in-time-compute model.

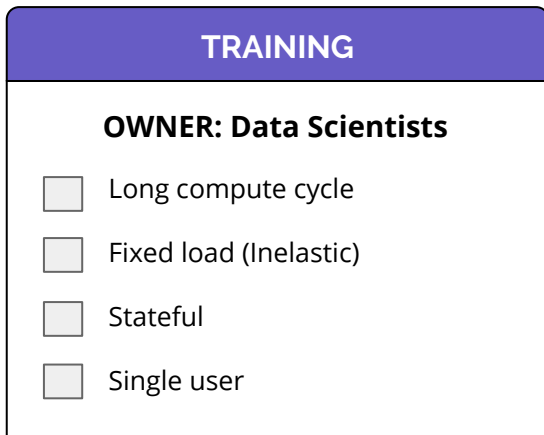
ADVANTAGES

- Cost / Efficiency
- Concurrency built-in
- Speed of development
- Improved latency

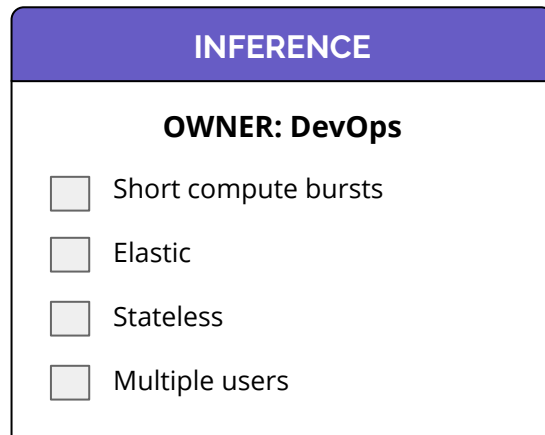


Analogous to dev tool chain.

Building and iterating over a model is similar to building an app.



Analogous to dev tool chain.
Building and iterating over a model
is similar to building an app.

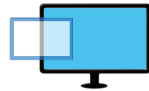


Analogous to an OS.
Running concurrent models
requires task scheduling.

TRAINING

OWNER: Data Scientists

- ☐ Long compute cycle
- ☐ Fixed load (Inelastic)
- ☐ Stateful
- ☐ Single user



Metal or VM

INFERENCE

OWNER: DevOps

- ☒ Short compute bursts
- ☐ Elastic
- ☒ Stateless
- ☐ Multiple users



Containers

TRAINING

OWNER: Data Scientists

- ☐ Long compute cycle
- ☐ Fixed load (Inelastic)
- ☐ Stateful
- ☐ Single user

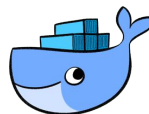


Metal or VM

INFERENCE

OWNER: DevOps

- ☒ Short compute bursts
- ☒ Elastic
- ☒ Stateless
- ☐ Multiple users



Containers

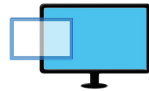


Kubernetes

TRAINING

OWNER: Data Scientists

- ☐ Long compute cycle
- ☐ Fixed load (Inelastic)
- ☐ Stateful
- ☐ Single user

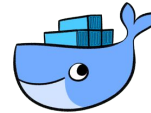


Metal or VM

INFERENCE

OWNER: DevOps

- ☒ Short compute bursts
- ☒ Elastic
- ☒ Stateless
- ☒ Multiple users



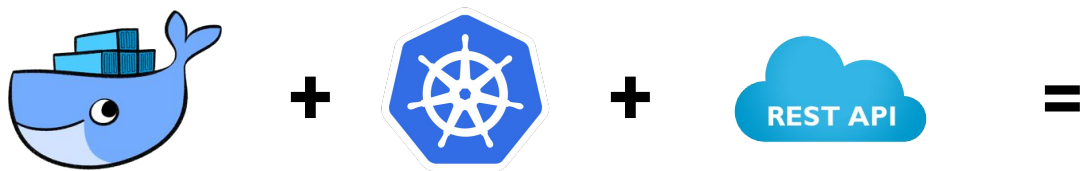
Containers



Kubernetes



Why Microservices?



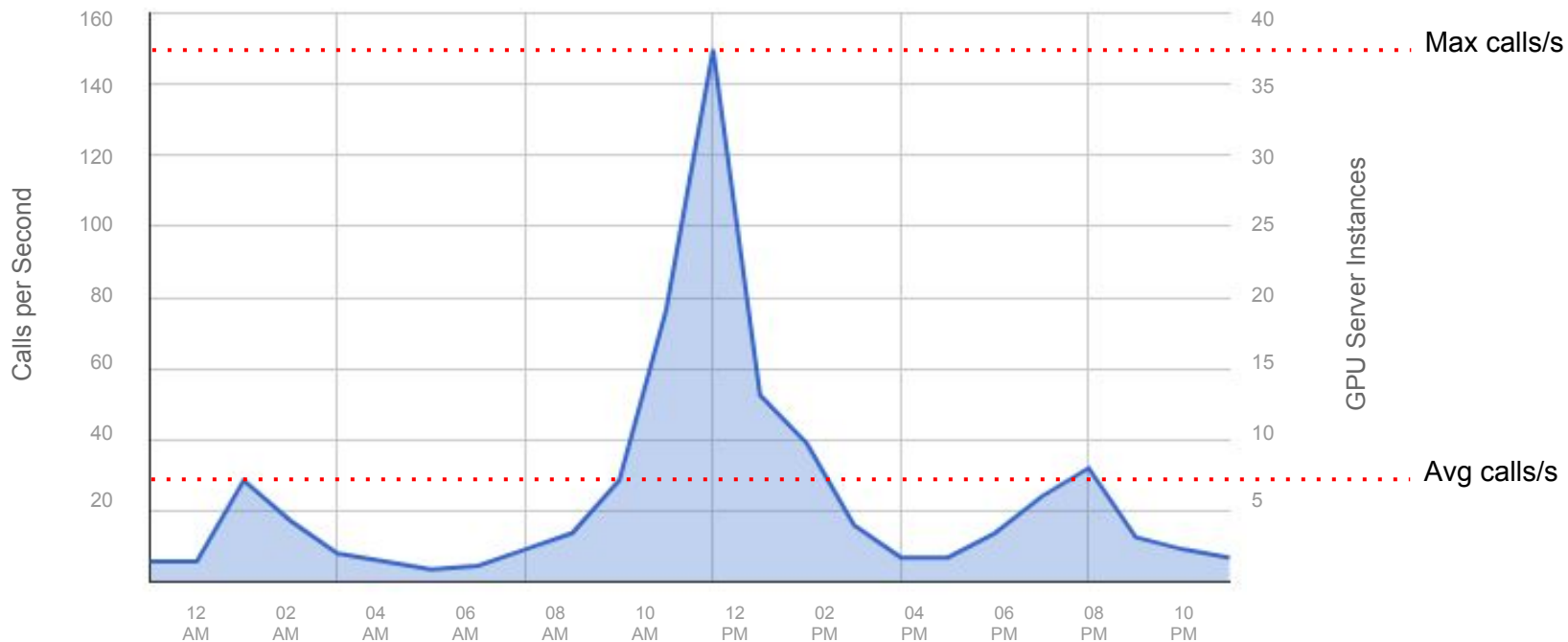
- Elastic
- Scalable
- Software agnostic
- Hardware agnostic

Why Serverless?

- Cost / Efficiency
- Concurrency built-in
- Improved latency

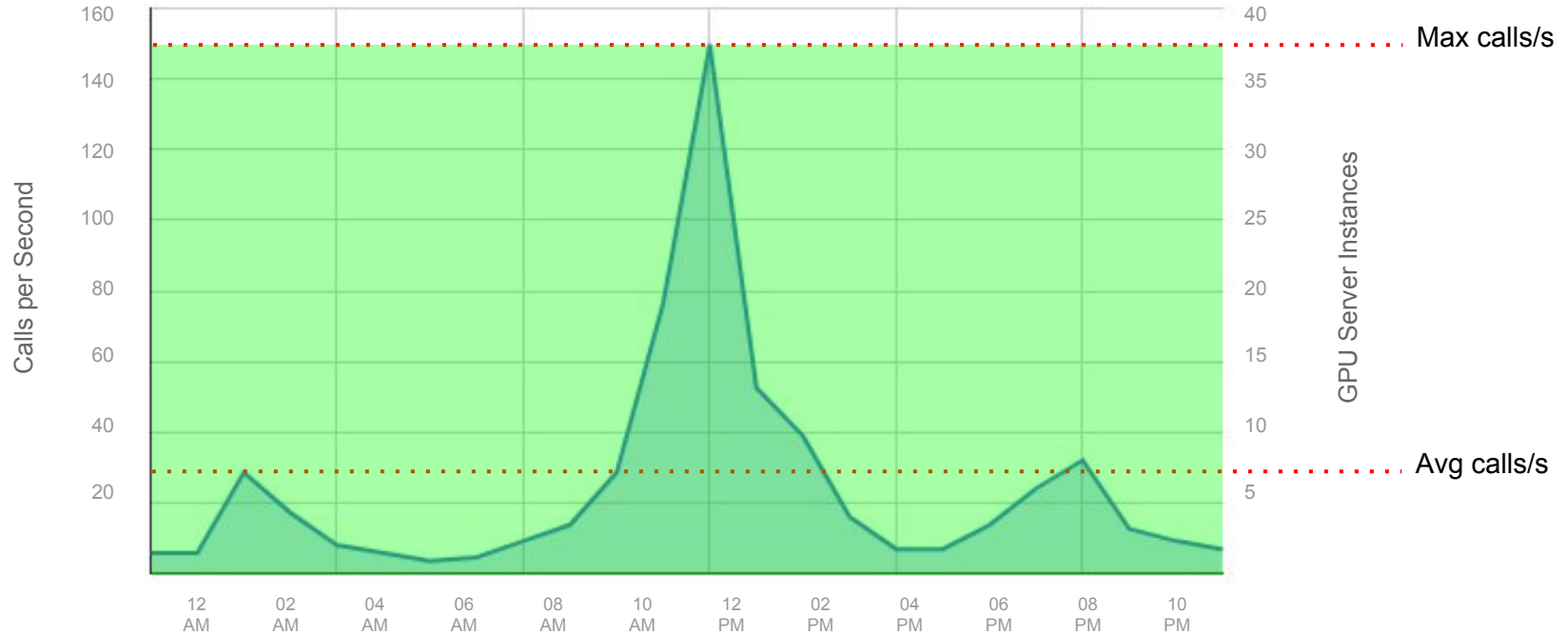
Why Serverless - Cost Efficiency

Jian Yang's "SeeFood" is most active during lunchtime.



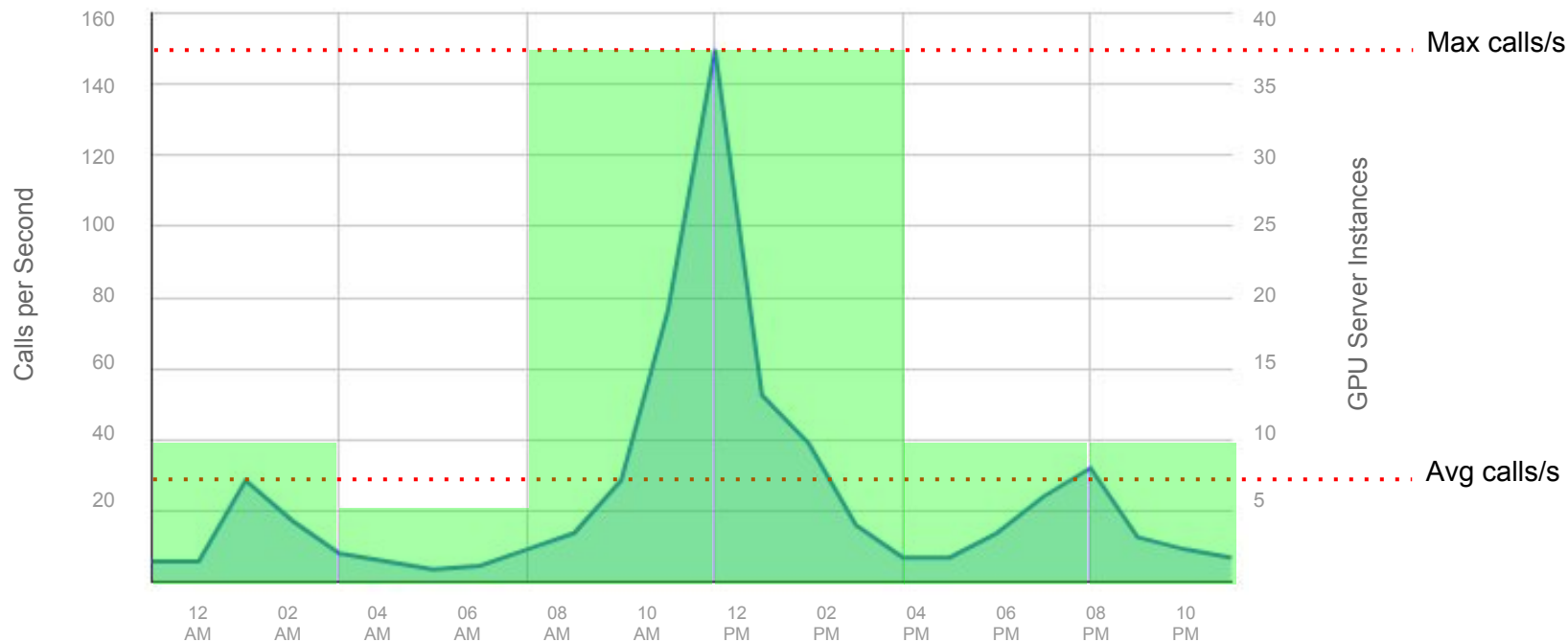
Traditional Architecture - Design for Maximum

40 machines 24 hours. $\$648 \times 40 = \$25,920$ per month



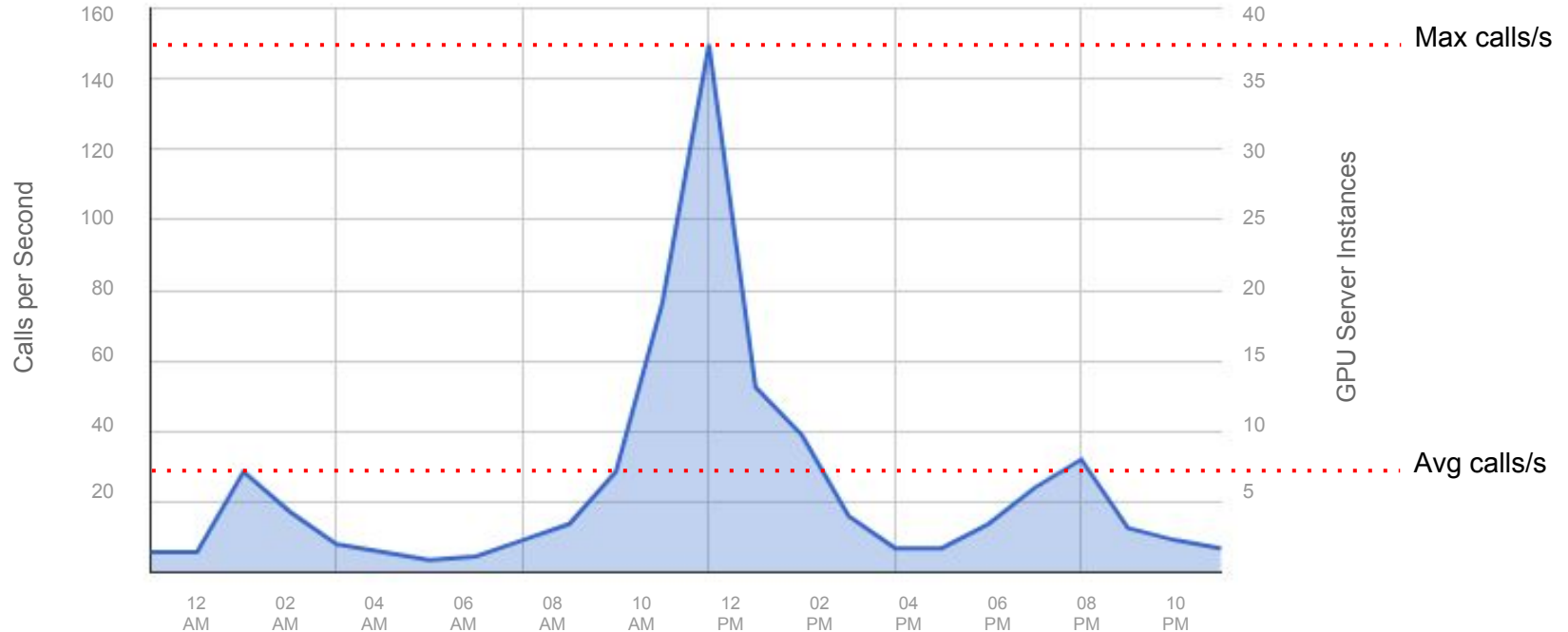
Autoscale Architecture - Design for Local Maximum

19 machines 24 hours. $\$648 \times 40 = \text{\$12,312 per month}$

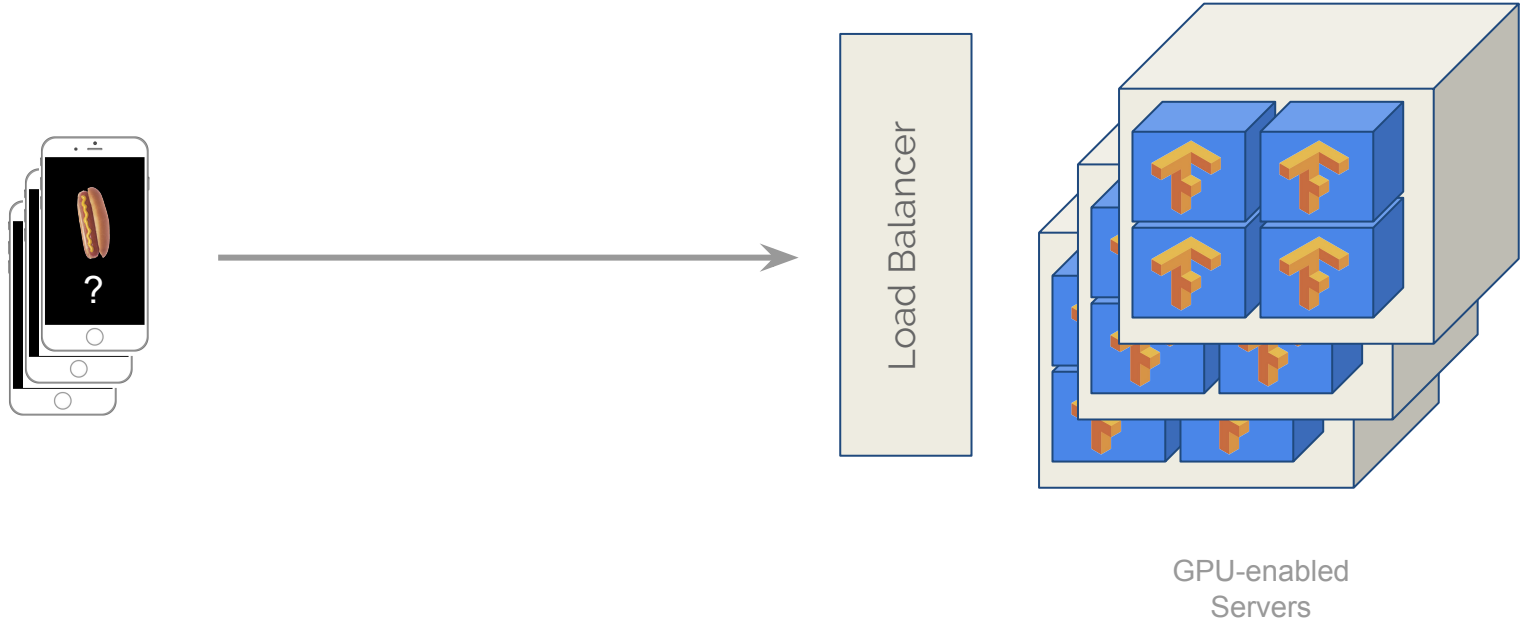


Serverless Architecture - Design for Minimum

Avg. of 21 calls / sec, or equivalent of 6 machines. $\$648 * 6 = \$3,888$ per month

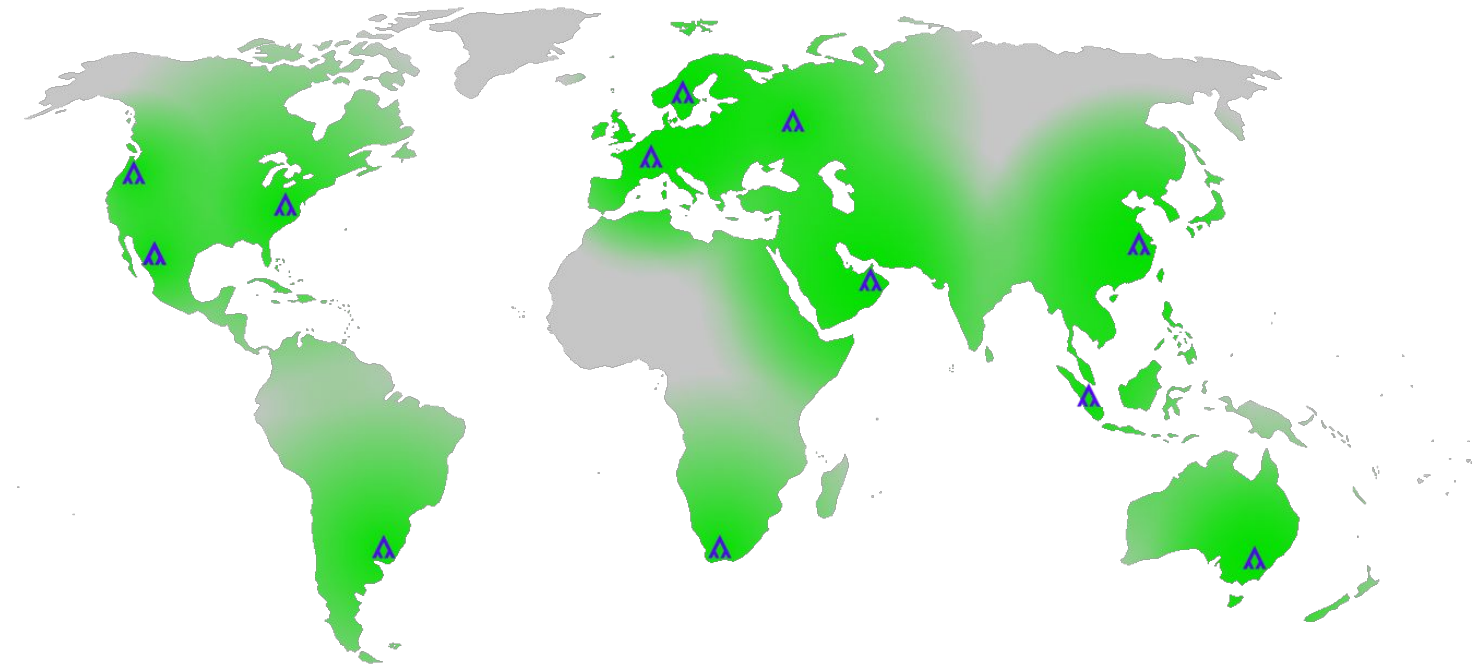


Why Serverless - Concurrency



Why Serverless - Improved Latency

Portability = Low Latency





ALSO:

GPU Memory Management, Job Scheduling, Cloud Abstraction, etc.

An Operating System for AI

op·er·at·ing sys·tem

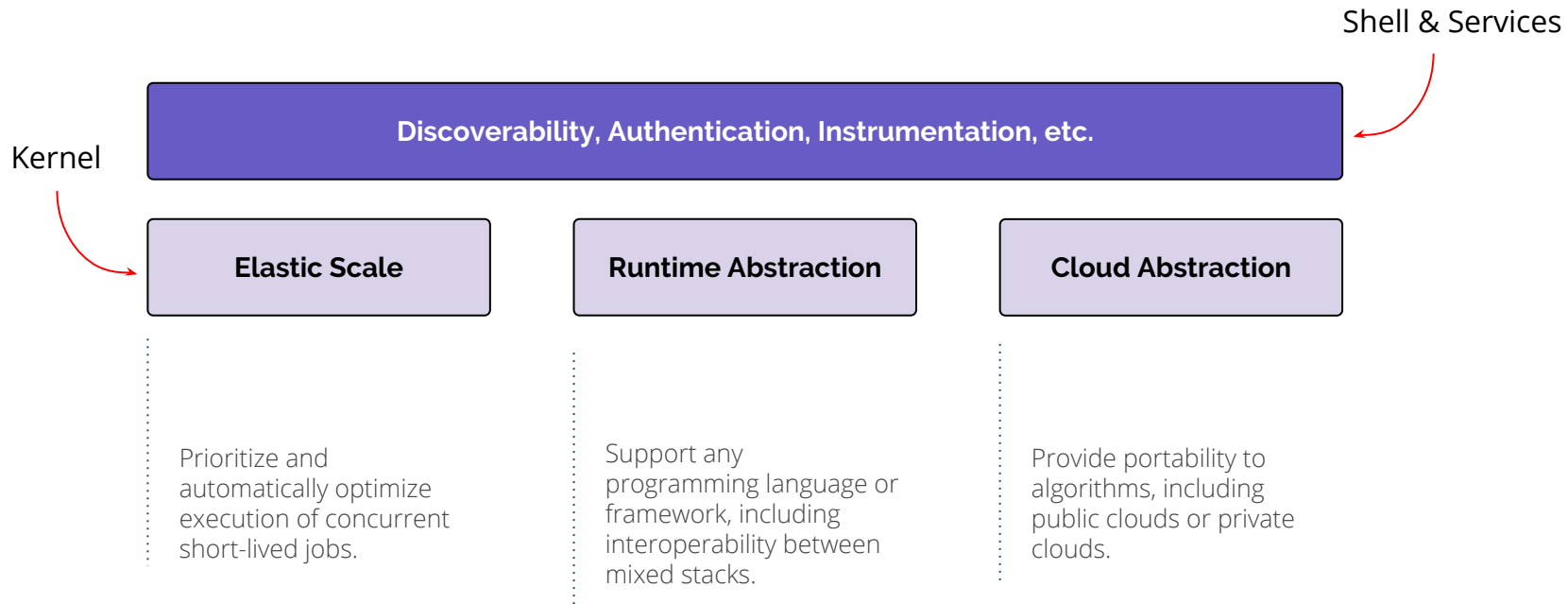
/ˈäpəˌrādiNG ˌsistəm/ 

noun

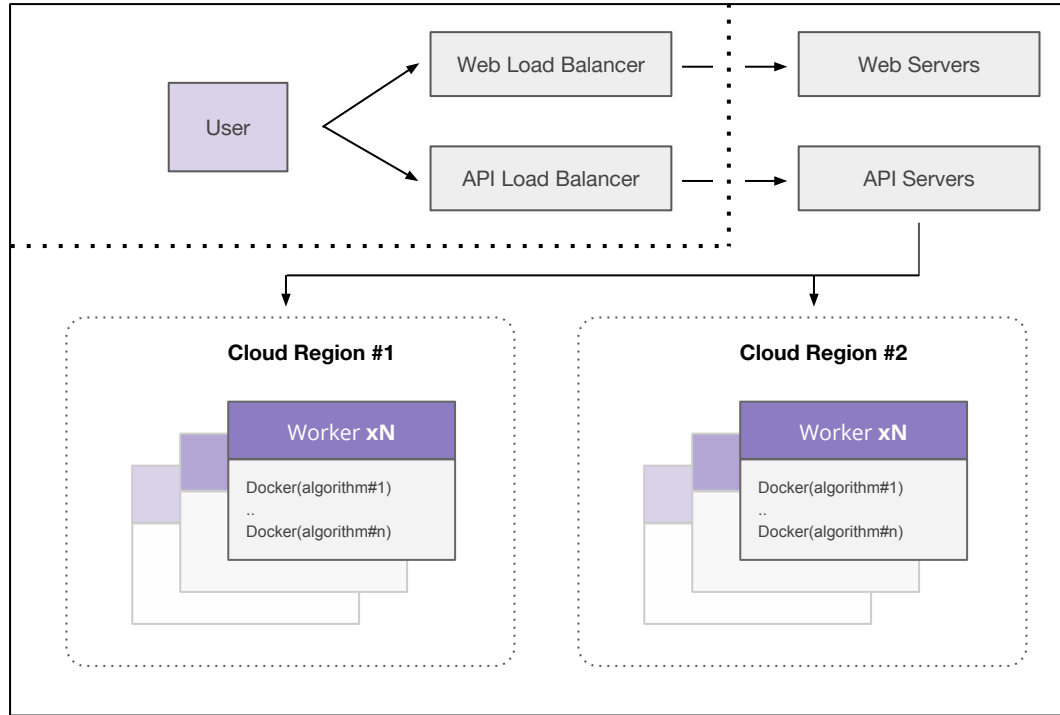
the software that supports a computer's basic functions, such as scheduling tasks, executing applications, and controlling peripherals.



Translations, word origin, and more definitions

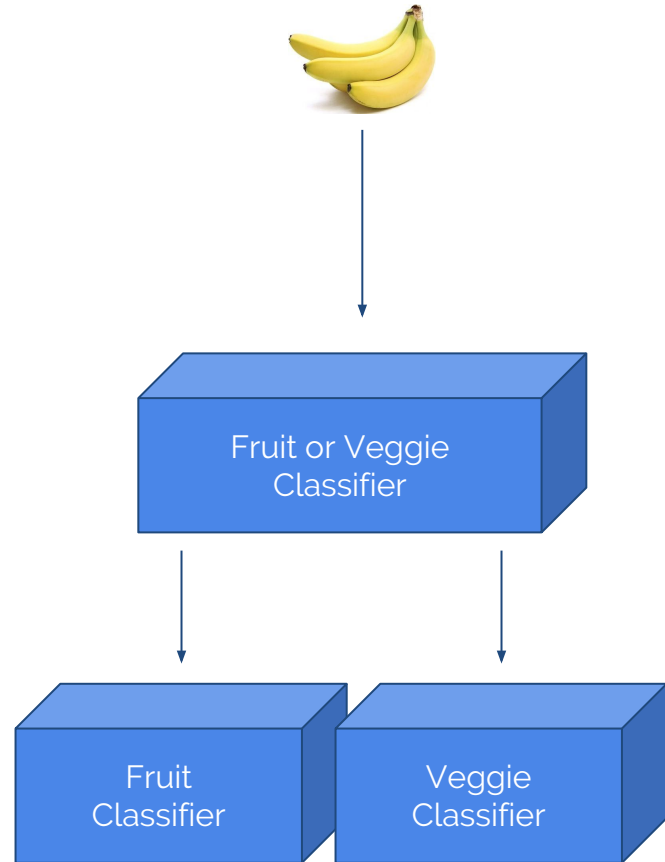


Kernel: Elastic Scale

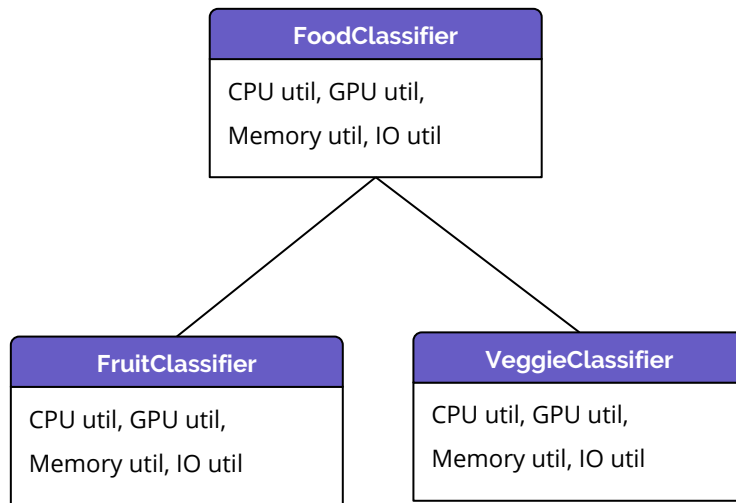


Composability

Composability is critical for AI workflows because of data processing pipelines and ensembles.



Kernel: Elastic Scale + Intelligent Orchestration



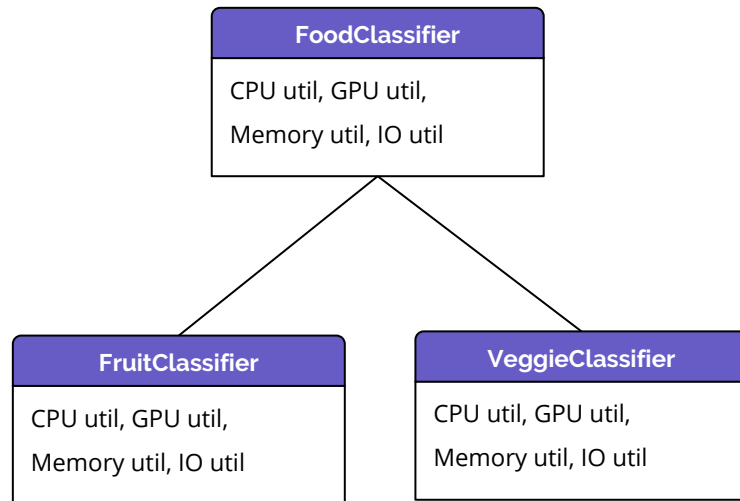
Kernel: Elastic Scale + Intelligent Orchestration

Knowing that:

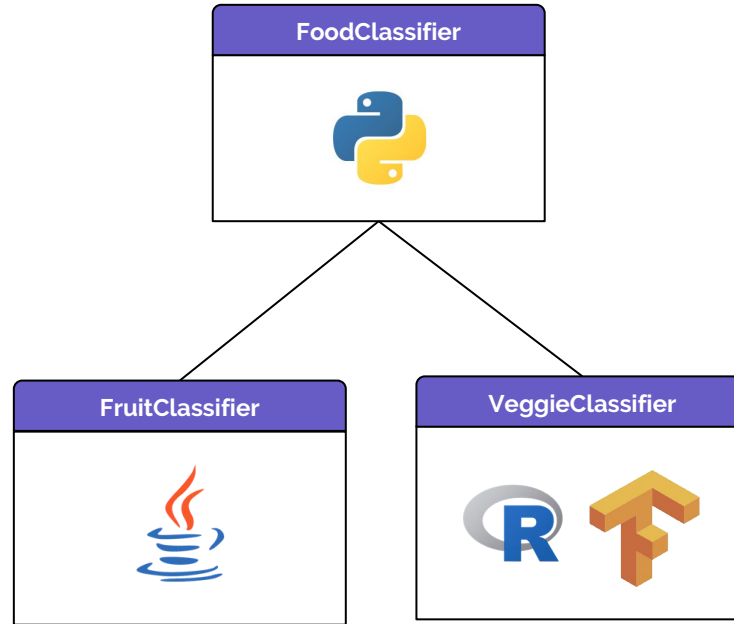
- **Algorithm A** always calls **Algorithm B**
- Algorithm A consumes X CPU, X Memory, etc
- Algorithm B consumes X CPU, X Memory, etc

Therefore we can slot them in a way that:

- Reduce network latency
- Increase cluster utilization
- Build dependency graphs



Kernel: Runtime Abstraction



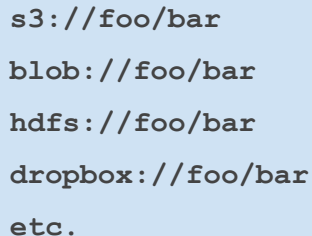
Kernel: Cloud Abstraction - Storage

No storage abstraction

```
s3      = boto3.client( "s3" )  
obj     = s3.get_object( Bucket= "bucket-name" , Key= "records.csv" )  
data    = obj[ "Body" ].read()
```





With storage abstraction

```
data    = Algorithmia().client.file( "blob://records.csv" ).get()
```



A diagram illustrating storage abstraction. A dashed vertical line connects the underlined path `blob://records.csv` from the code block above to a light blue rectangular box. Inside the box, a list of storage paths is shown, all starting with `//foo/bar`: `s3://foo/bar`, `blob://foo/bar`, `hdfs://foo/bar`, `dropbox://foo/bar`, and `etc.`

Kernel: Cloud Abstraction

		 Google Cloud Platform		
Compute	EC2	CE	VM	Nova
Autoscaling	Autoscaling Group	Autoscaler	Scale Set	Heat Scaling Policy
Load Balancing	Elastic Load Balancer	Load Balancer	Load Balancer	LBaaS
Remote Storage	Elastic Block Store	Persistent Disk	File Storage	Block Storage

Summary - What makes an OS for AI?

Stack-agnostic

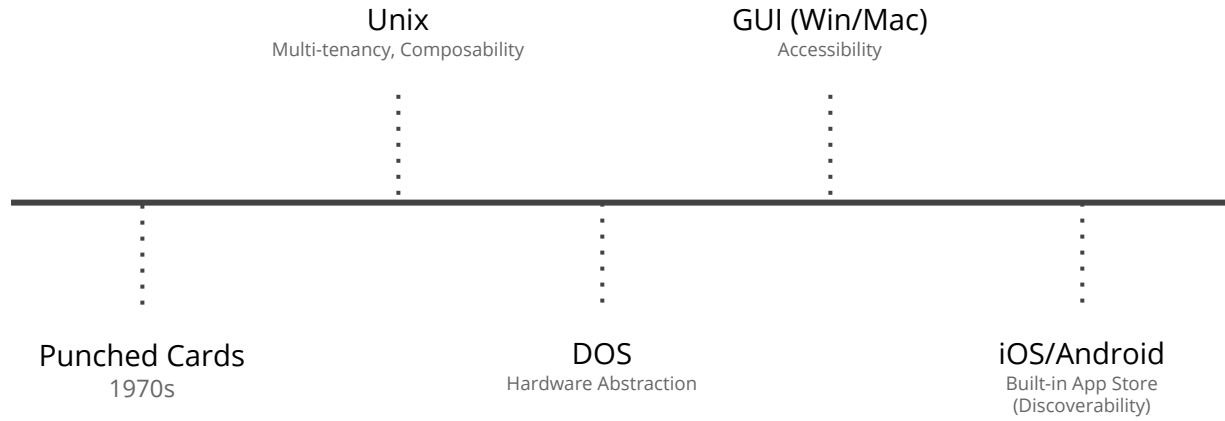
Composable

Self-optimizing

Auto-scaling

Monitorable

Discoverability



AI is here



FREE STUFF:

Signup with code:
CloudSummit17 for \$50 on us.

ALGORITHMIA

Thank you!

Diego Oppenheimer
CEO



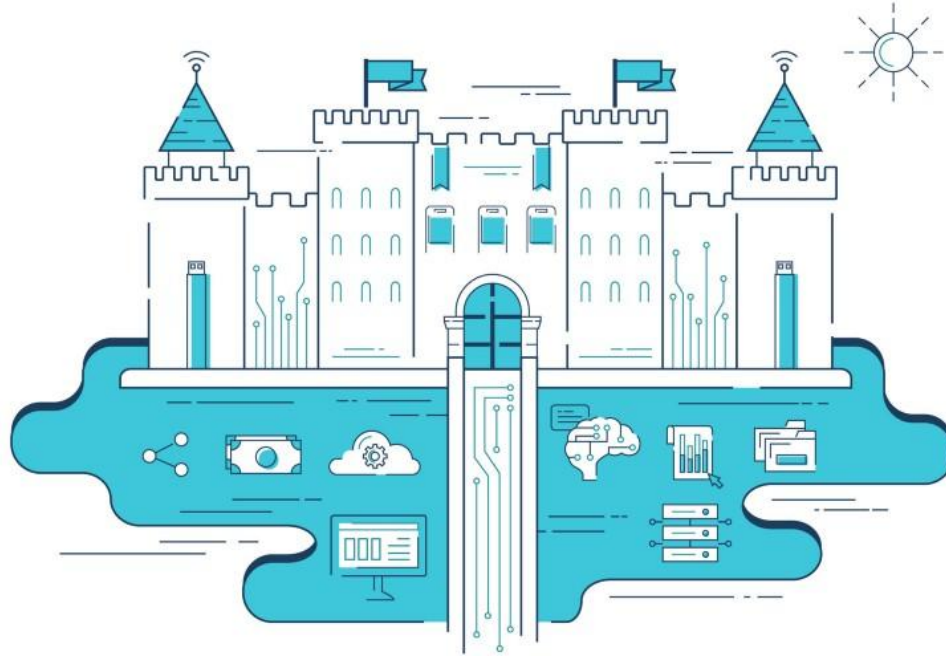
diego@algorithmia.com



@doppenhe

more slides

The New Moats





Kernel: Cloud Abstraction - Storage

```
# init
client = Algorithmia.client()

# get data (S3)
s3 = boto3.client( "s3" )
obj = s3.get_object( Bucket= "bucket-name" ,
Key="records.csv" )
data = obj[ "Body" ].read()

# remove seasonality
data = client.algo( "ts/RemoveSeasonality" ).pipe(data).result

# forecast time series
data = client.algo( "ts/ForecastLSTM" ).pipe(data).result
```

```
# init
client = Algorithmia.client()

# get data (anything)
data = client.file( "blob://records.csv" ).get()

# remove seasonality
data = client.algo( "ts/RemoveSeasonality" ).pipe(data).result

# forecast time series
data = client.algo( "ts/ForecastLSTM" ).pipe(data).result
```

```
01  # MY_ALGORITHM.py
02
03  client = Algorithmia.client()
04  data   = client.file("blob://records.csv").get()
05
06  # remove seasonality
07  data = client.algo("ts/RemoveSeasonality").pipe(data).result
08
09  # forecast time series
10  data = client.algo("ts/ForecastLSTM").pipe(data).result
```

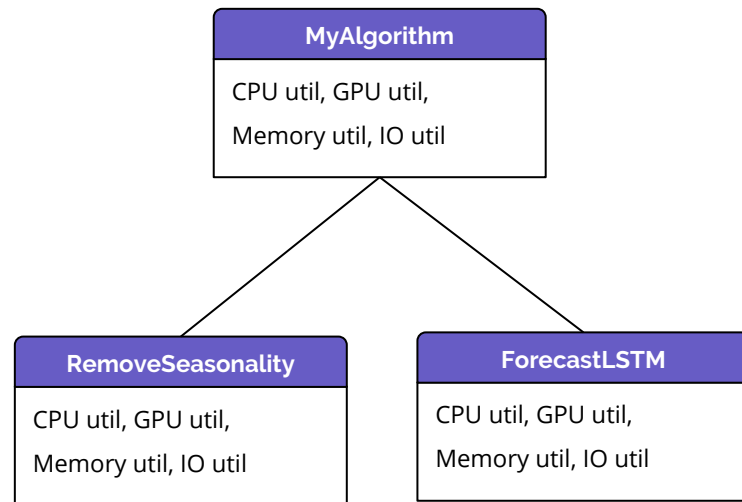
Kernel: Elastic Scale + Intelligent Orchestration

```
# MY_ALGORITHM.py

client = Algorithmia.client()
data = client.file( "blob://records.csv" ).get()

# remove seasonality
data = client.algo( "ts/RemoveSeasonality" ).pipe(data).result

# forecast time series
data = client.algo( "ts/ForecastLSTM" ).pipe(data).result
```



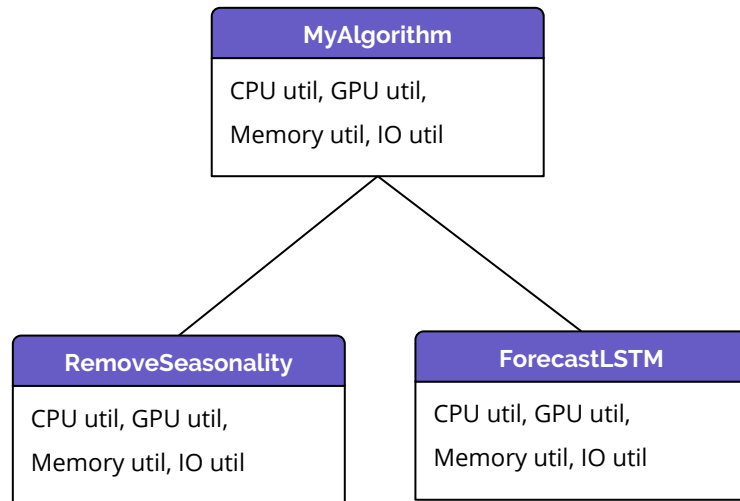
Kernel: Elastic Scale + Intelligent Orchestration

Knowing that:

- **Algorithm A** always calls **Algorithm B**
- Algorithm A consumes X CPU, X Memory, etc
- Algorithm B consumes X CPU, X Memory, etc

Therefore we can slot them in a way that:

- Reduce network latency
- Increase cluster utilization
- Build dependency graphs



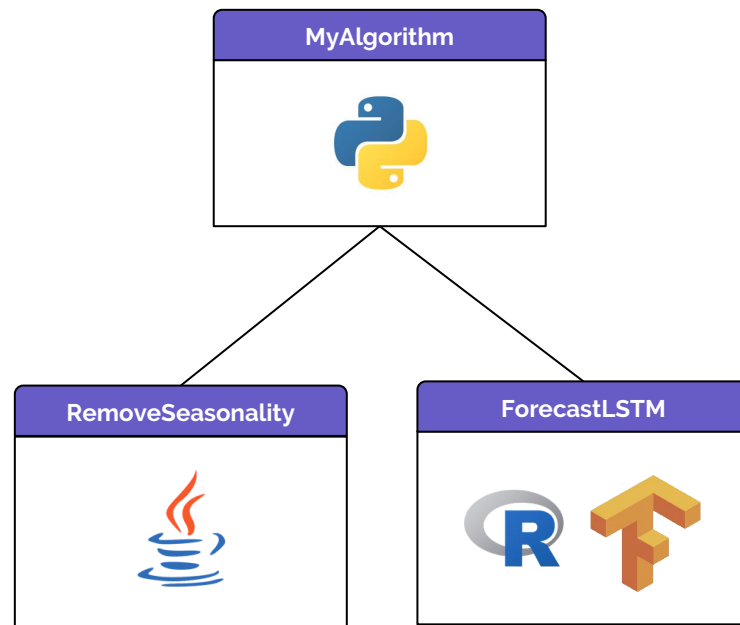
Kernel: Runtime Abstraction

```
# MY_ALGORITHM.py

client = Algorithmia.client()
data = client.file( "blob://records.csv" ).get()

# remove seasonality
data = client.algo( "ts/RemoveSeasonality" ).pipe(data).result

# forecast time series
data = client.algo( "ts/ForecastLSTM" ).pipe(data).result
```



Challenges

- **Machine learning**
 - CPU/GPU/Specialized hardware
 - Multiple frameworks, languages, dependencies
 - Called from different devices/architectures
- **“Snowflake” environments**
 - Unique cloud hardware and services
- **Uncharted territory**
 - Not a lot of literature, errors messages sometimes cryptic (can't just stackoverflow)