

runway

Building a Data Foundation for Multimodal Foundation Models

Ethan Rosenthal
Data Council 2025



About

Currently: ML @ Runway

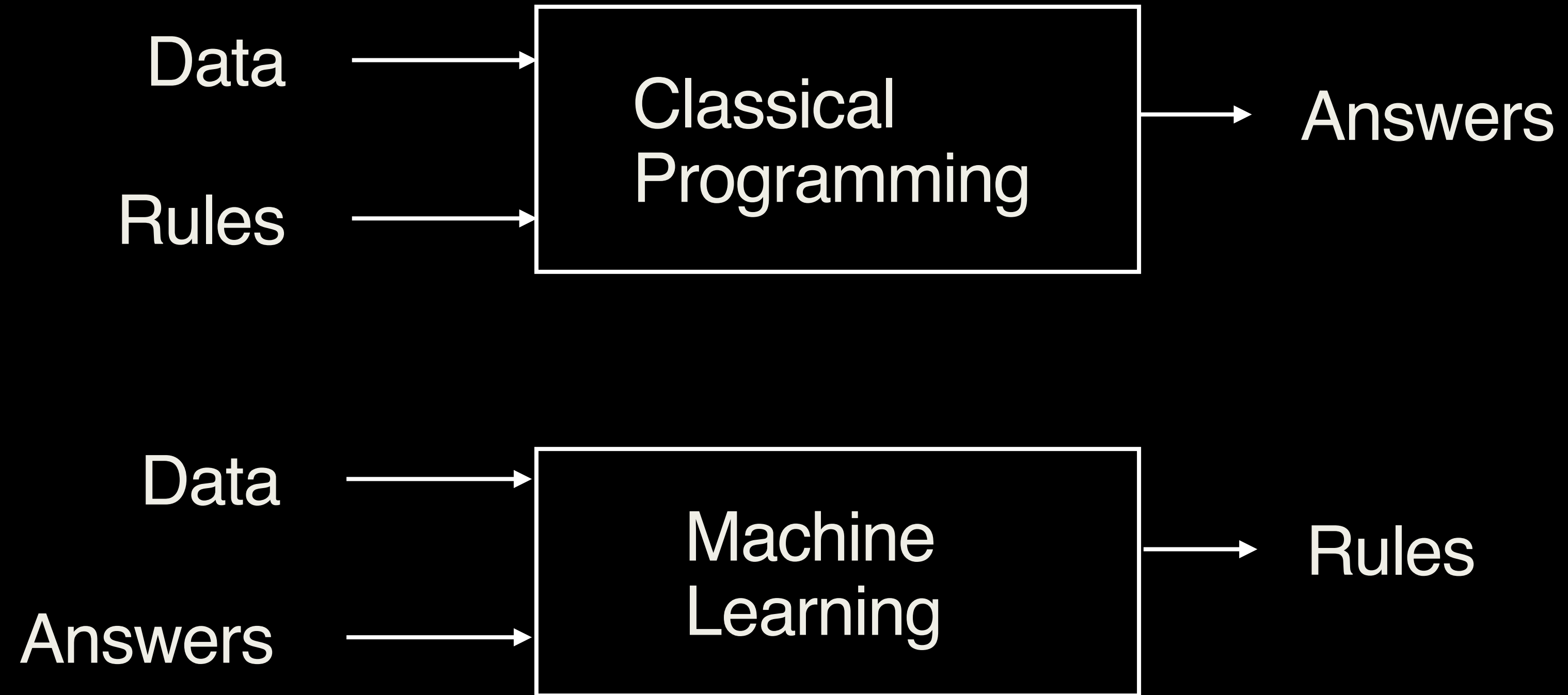
Formerly: LLMs @ Square; startups; physics PhD

Also: Adjunct Professor @ NYU Tandon

Find me



ethanrosenthal.com
[@ethanrosenthal.com](https://twitter.com/ethanrosenthal)
[@EthanRosenthal](https://github.com/EthanRosenthal)



Data



Rules



Classical
Programming



Answers

Data



Answers



Machine
Learning



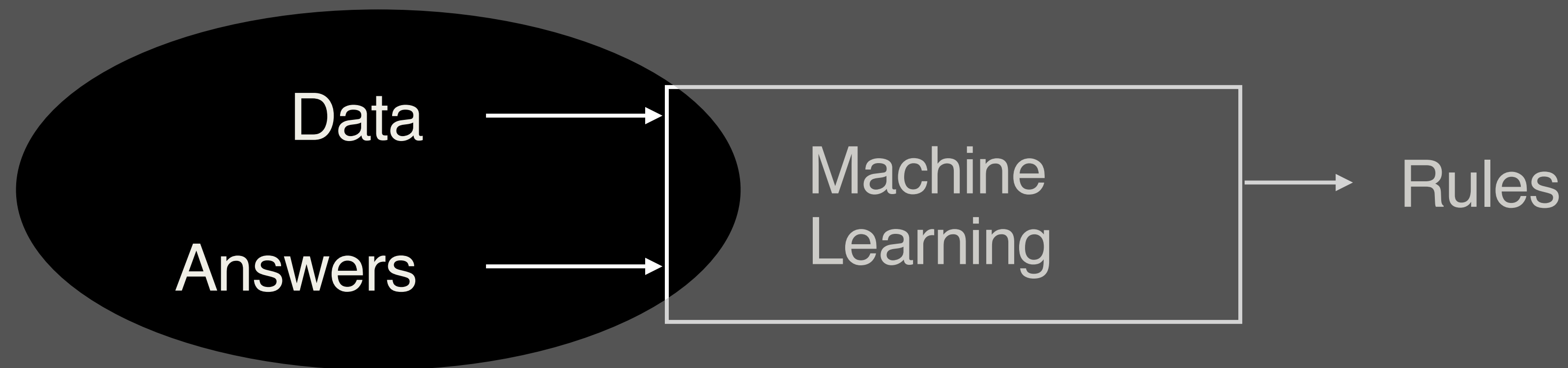
Rules

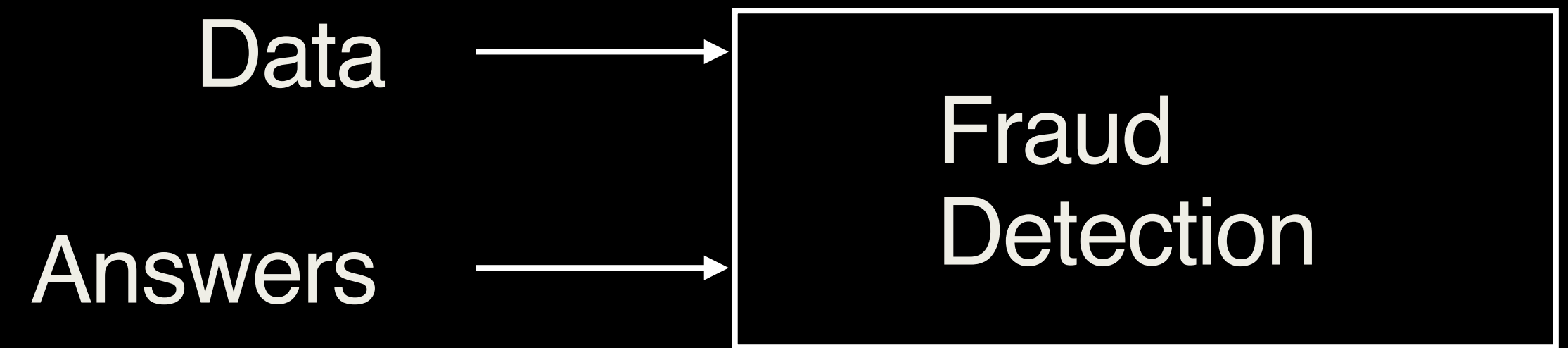
“Everyone wants to do the model work, not the data work”: **Data Cascades in High-Stakes AI**

Nithya Sambasivan, Shivani Kapania, Hannah Highfill, Diana Akrong, Praveen Paritosh, Lora Aroyo

[nithyasamba,kapania,hhighfill,dakrong,pkp,lora]@google.com

Google Research
Mountain View, CA

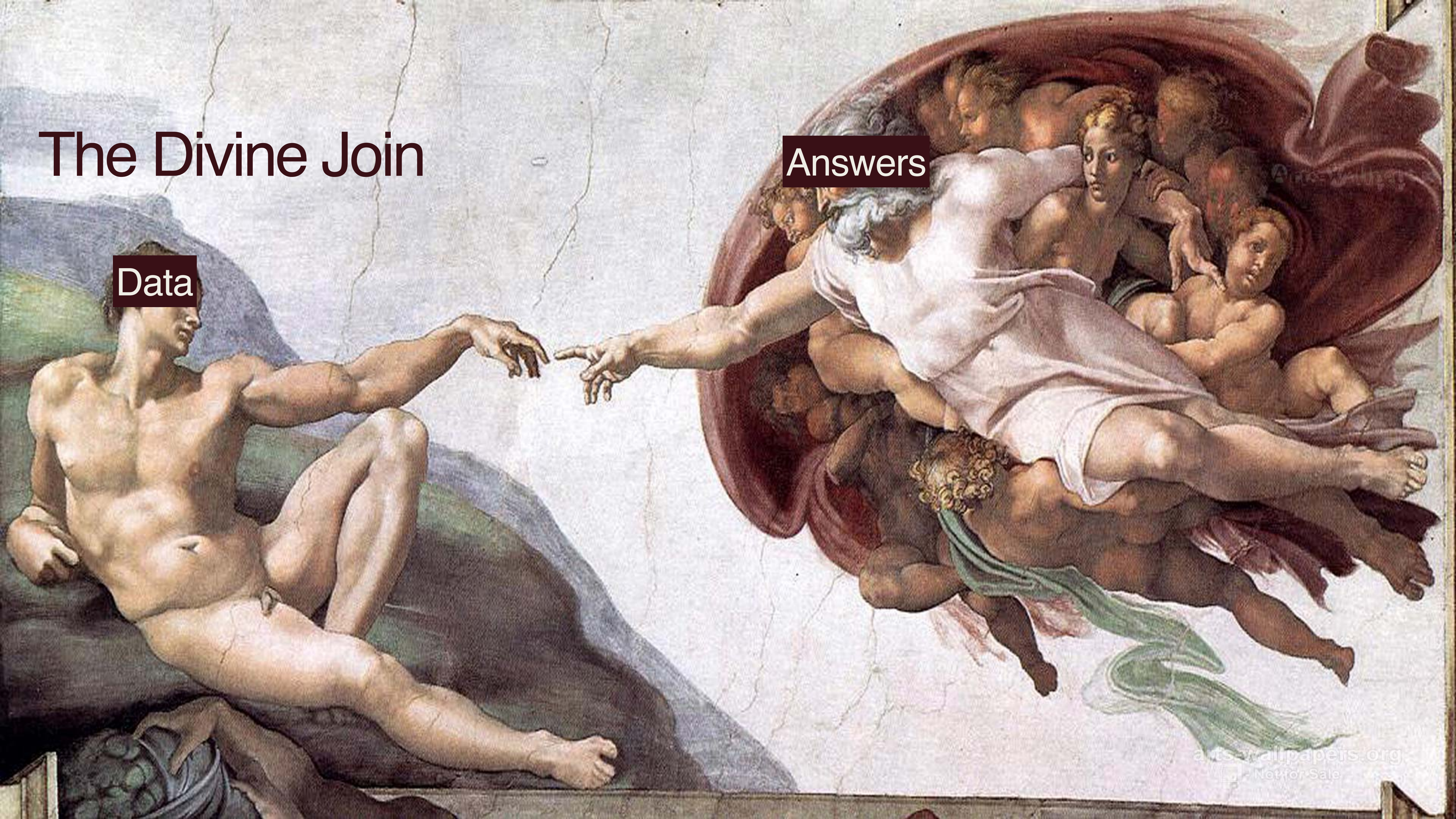




The Divine Join

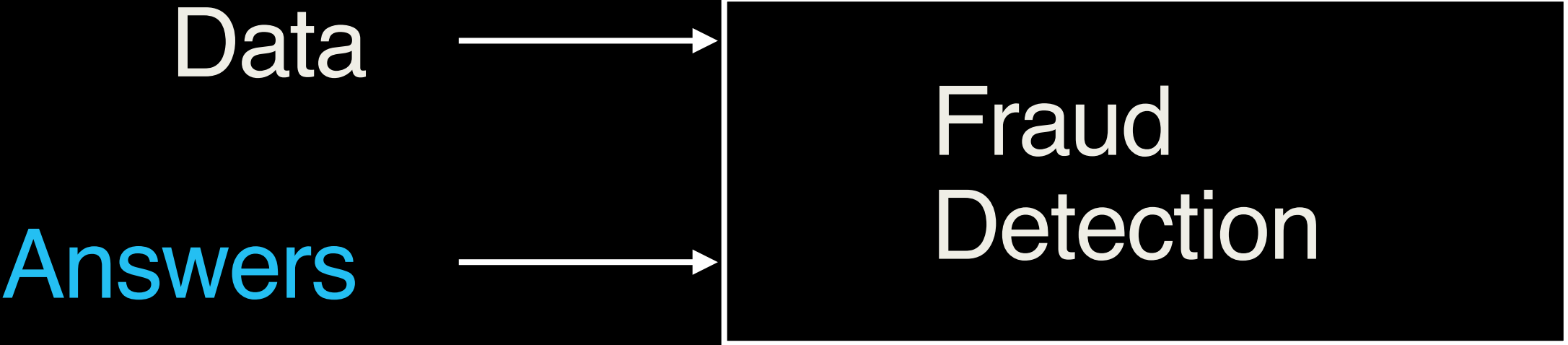
Data

Answers



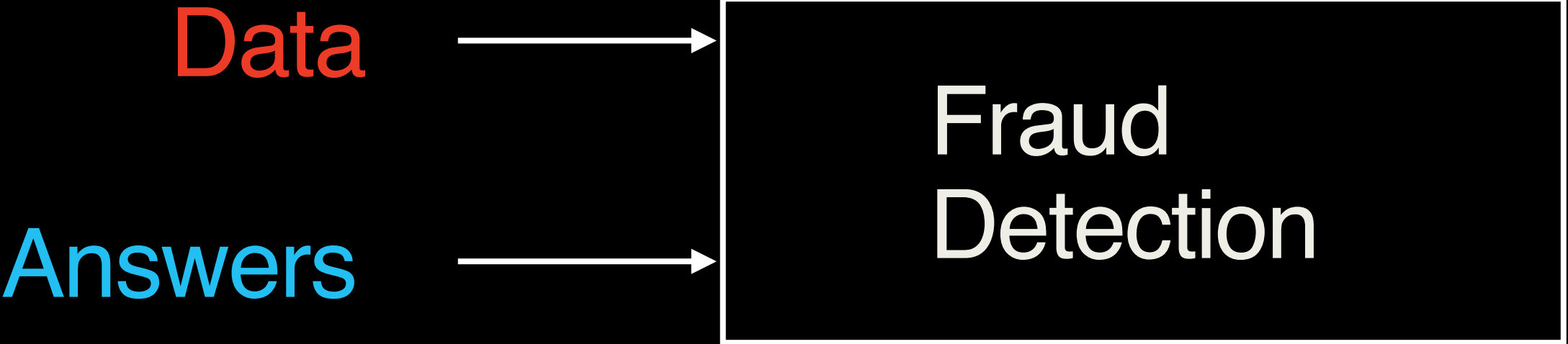
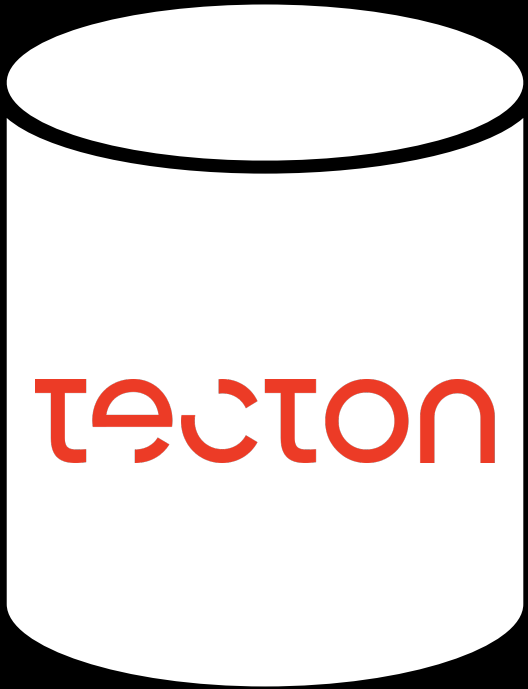


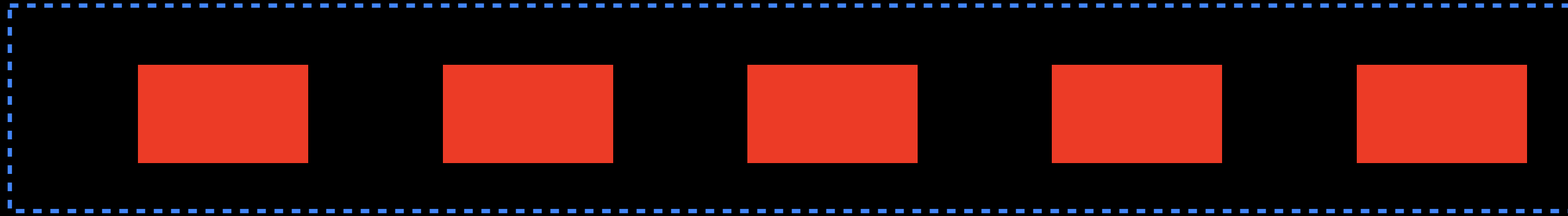
Is Fraud?

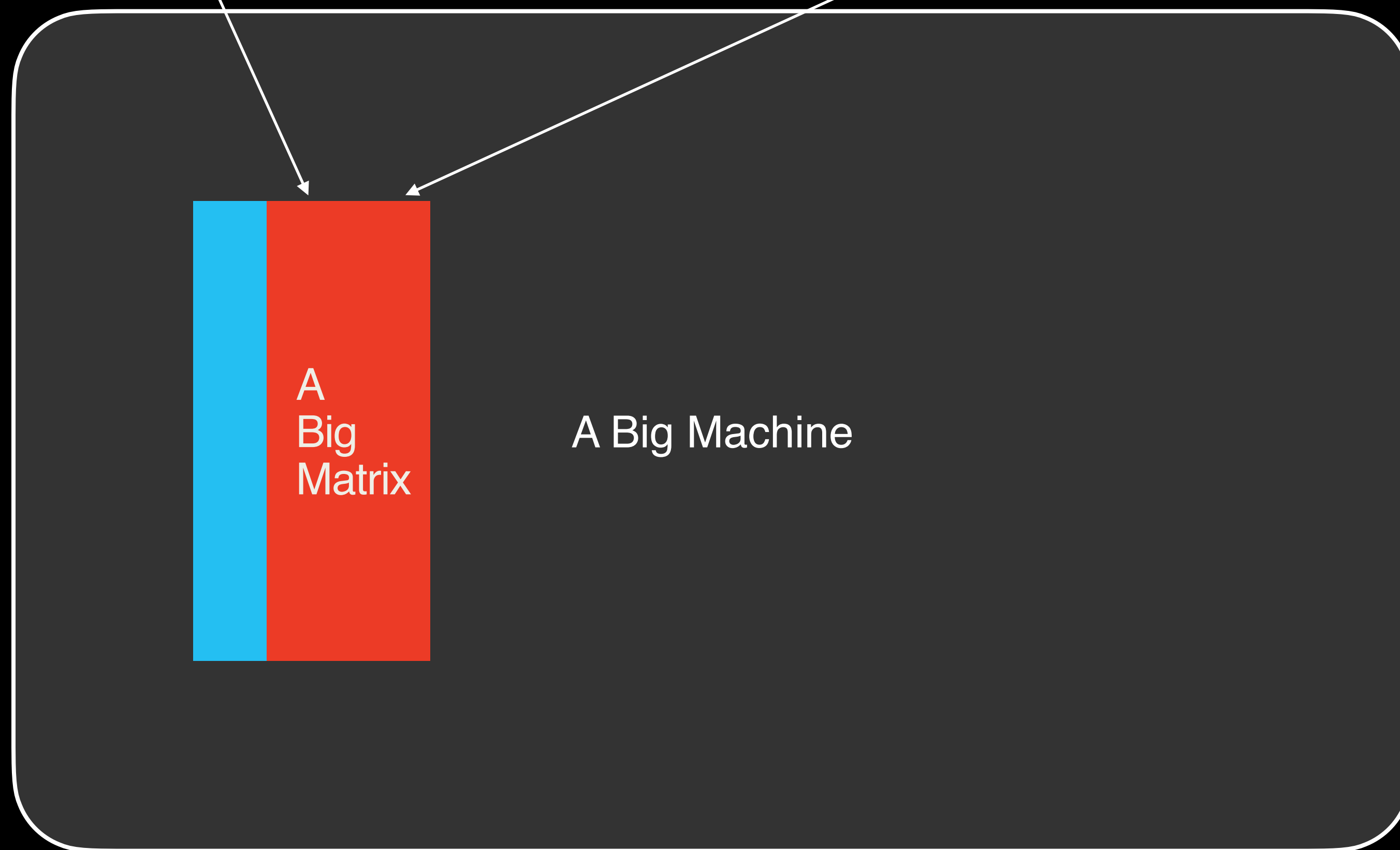
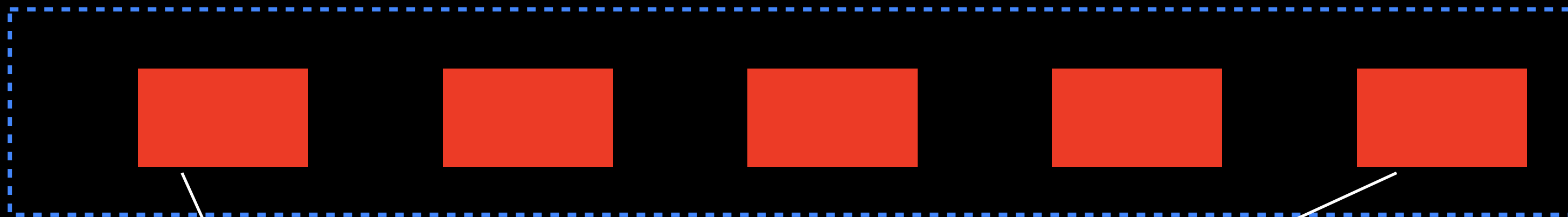


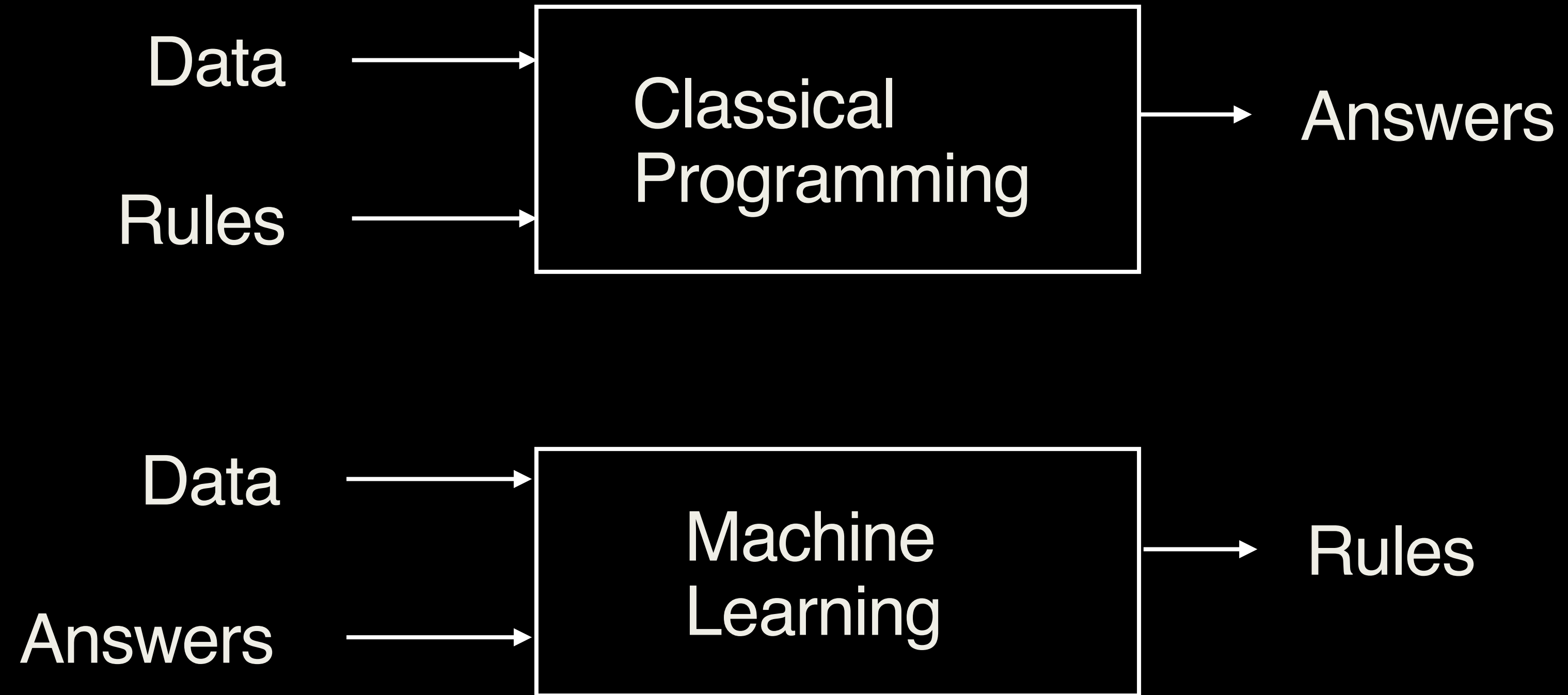


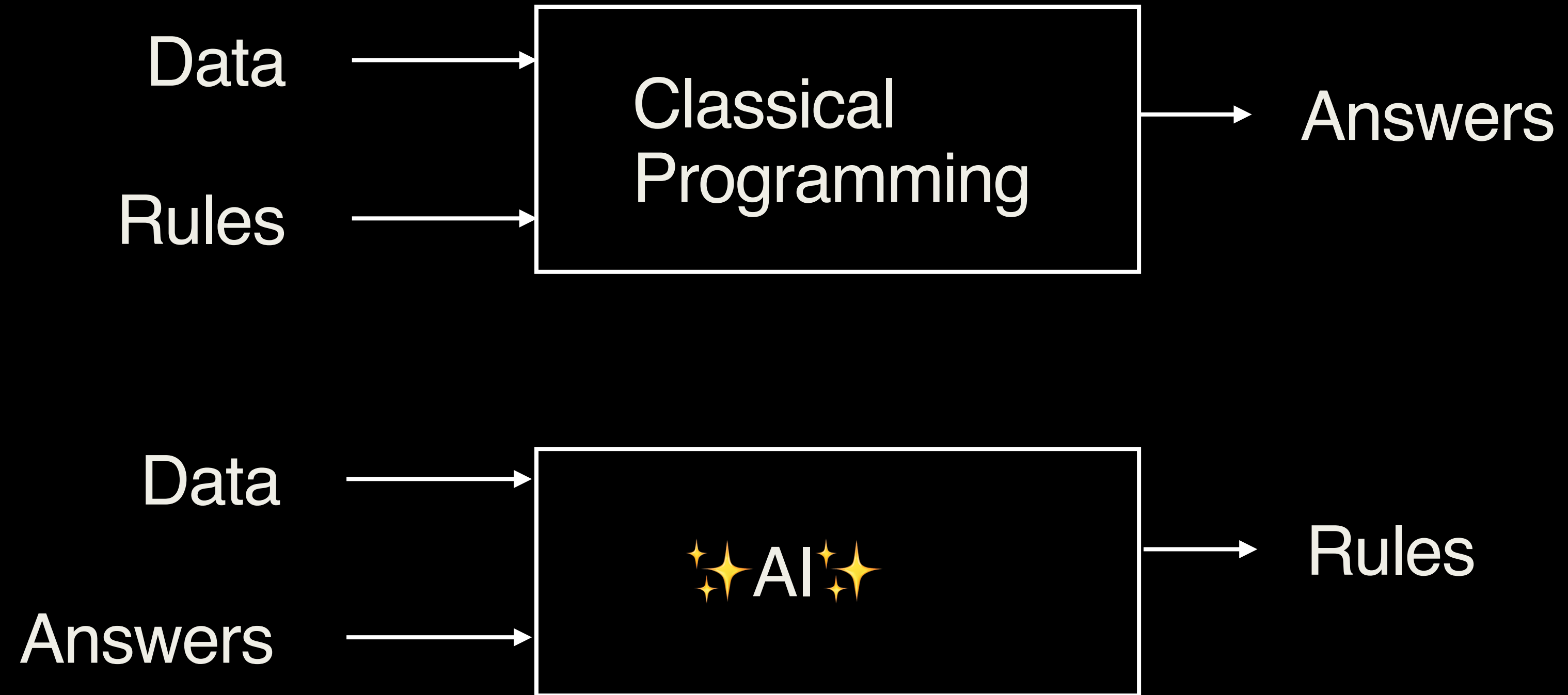
Is Fraud?	IP Address, qps, dollar amt, etc...			







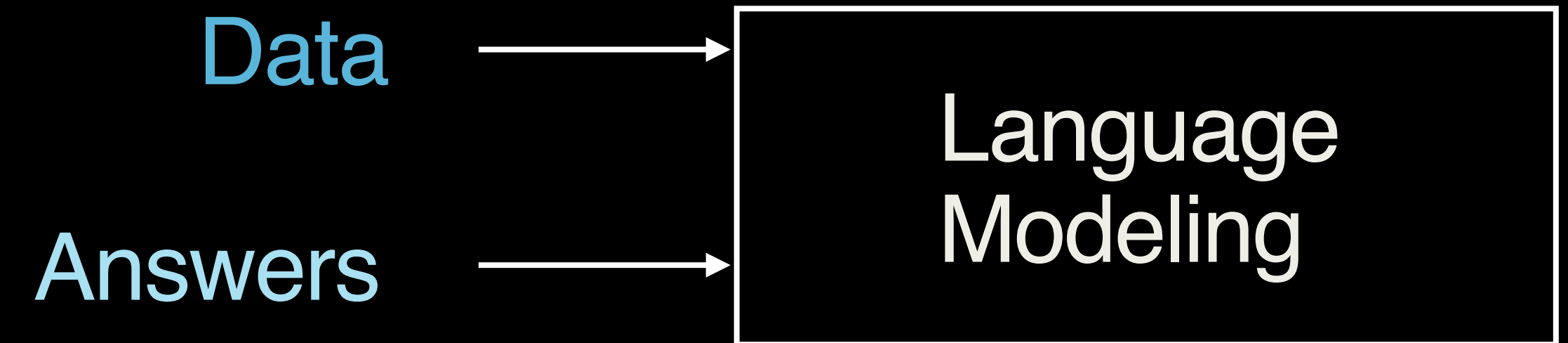


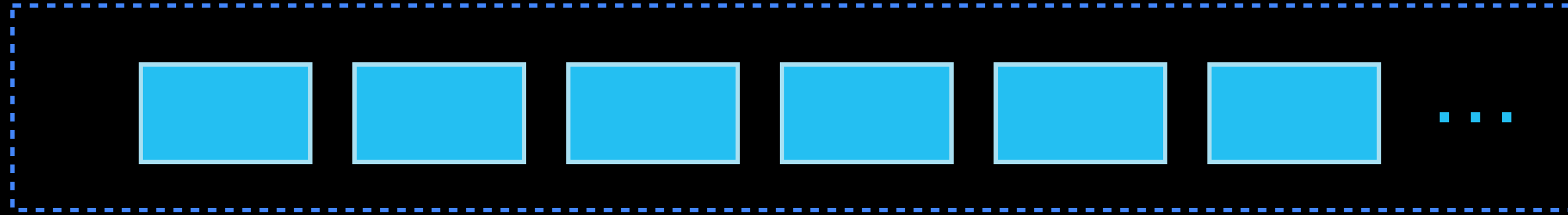


The Corpus Dump

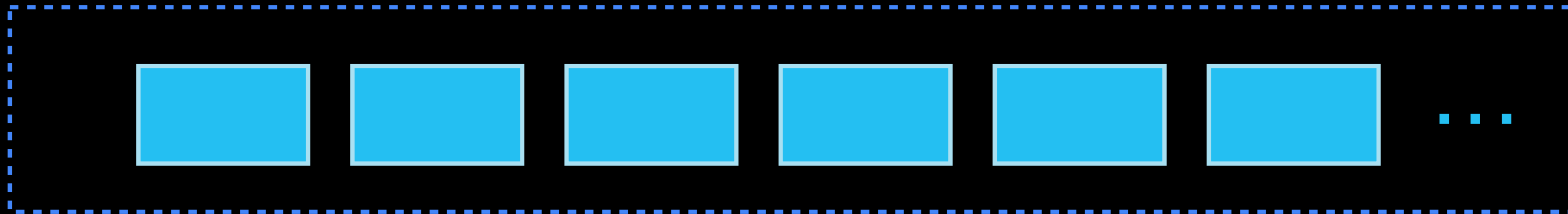


Customer: Do you all do kids birthdays?
Merchant: Of course!





A Big Machine



GPU Process 1

GPU Process 2

GPU Process 3

GPU Process 4

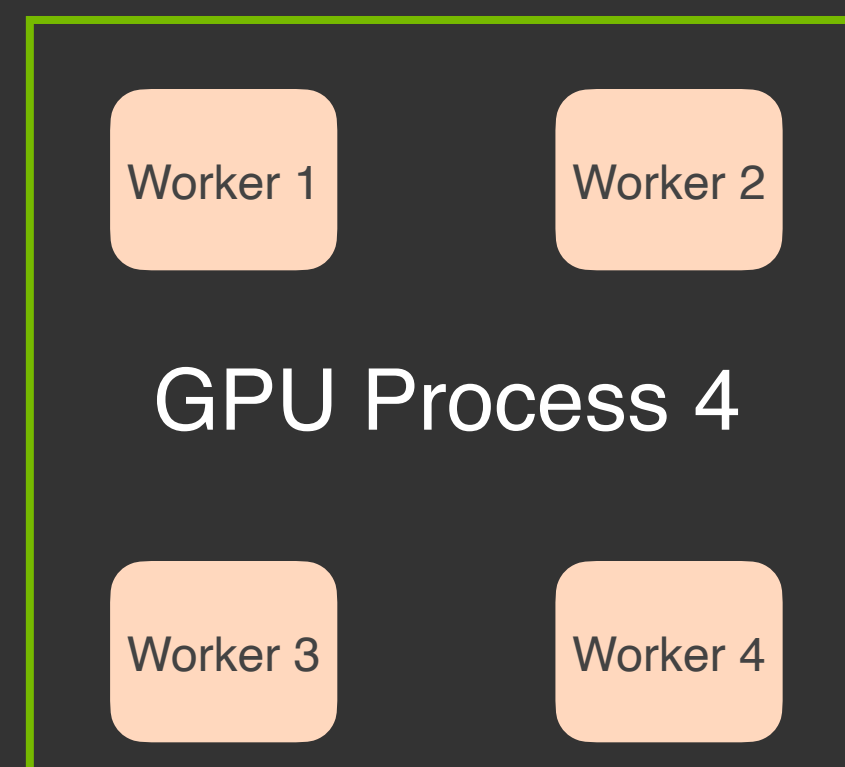
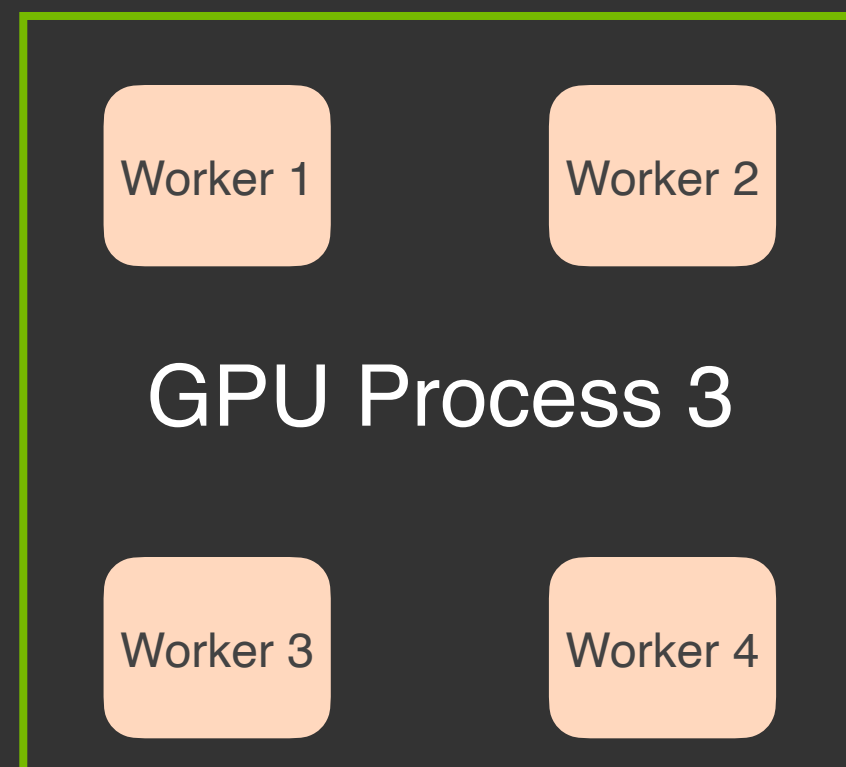
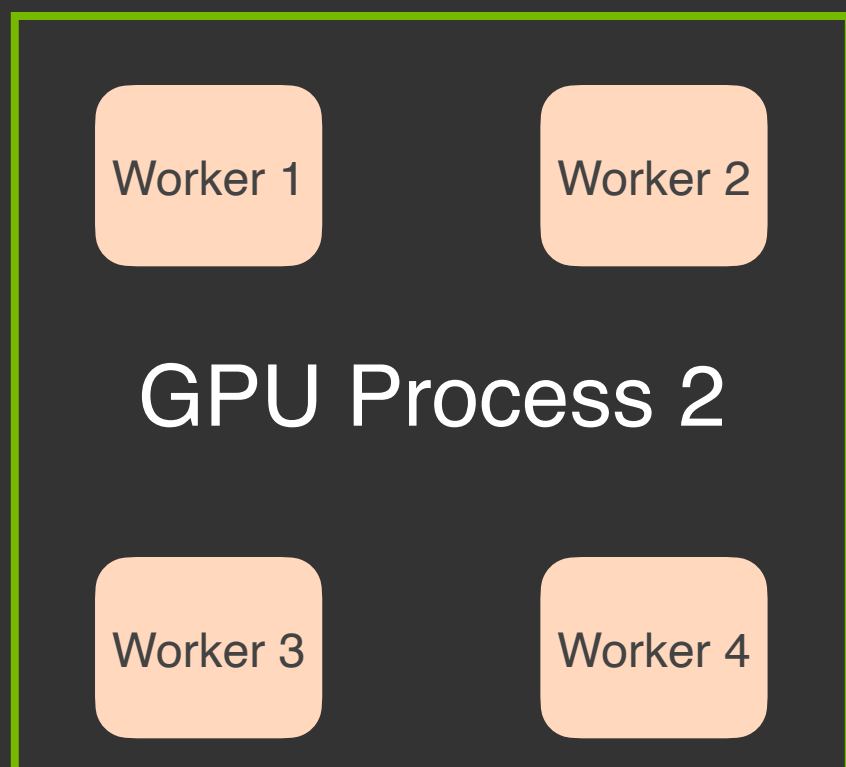
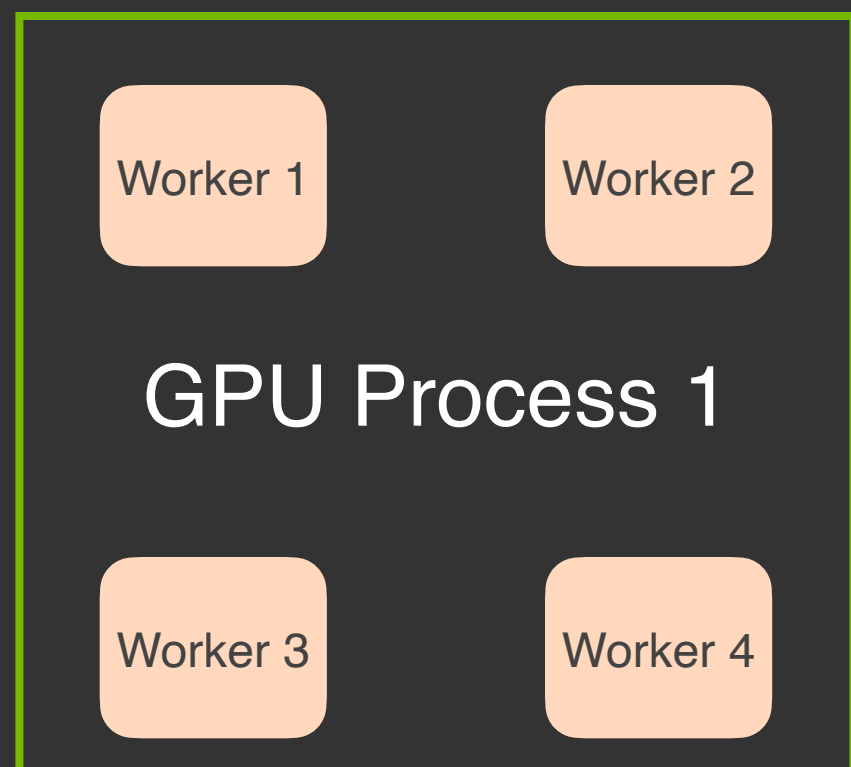
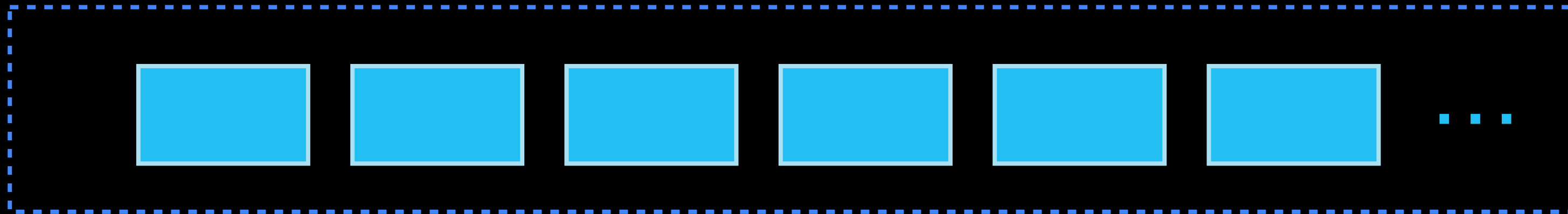
A Big Machine

GPU Process 5

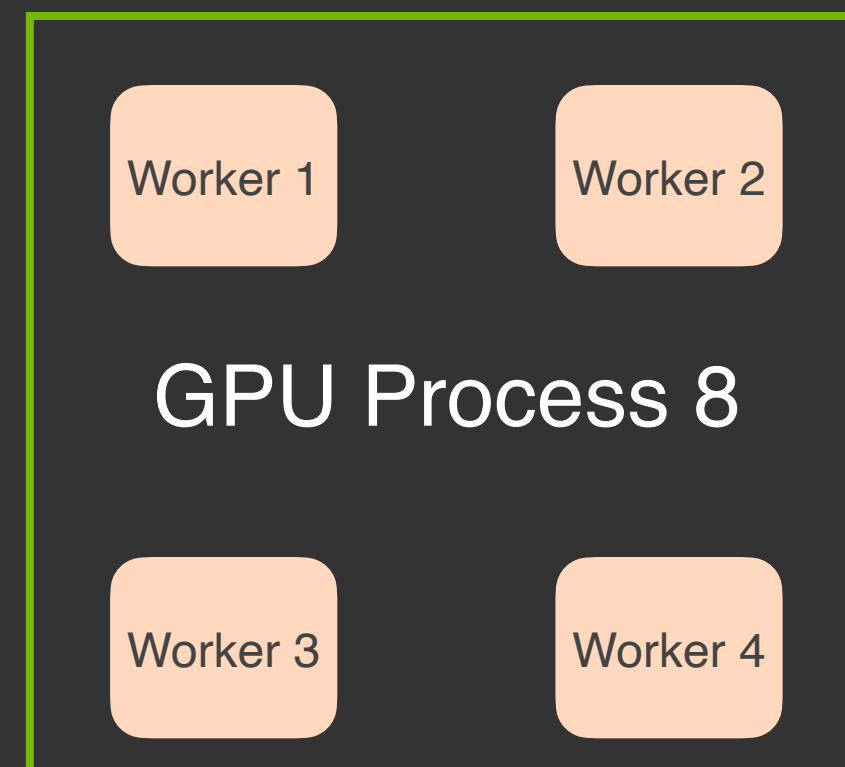
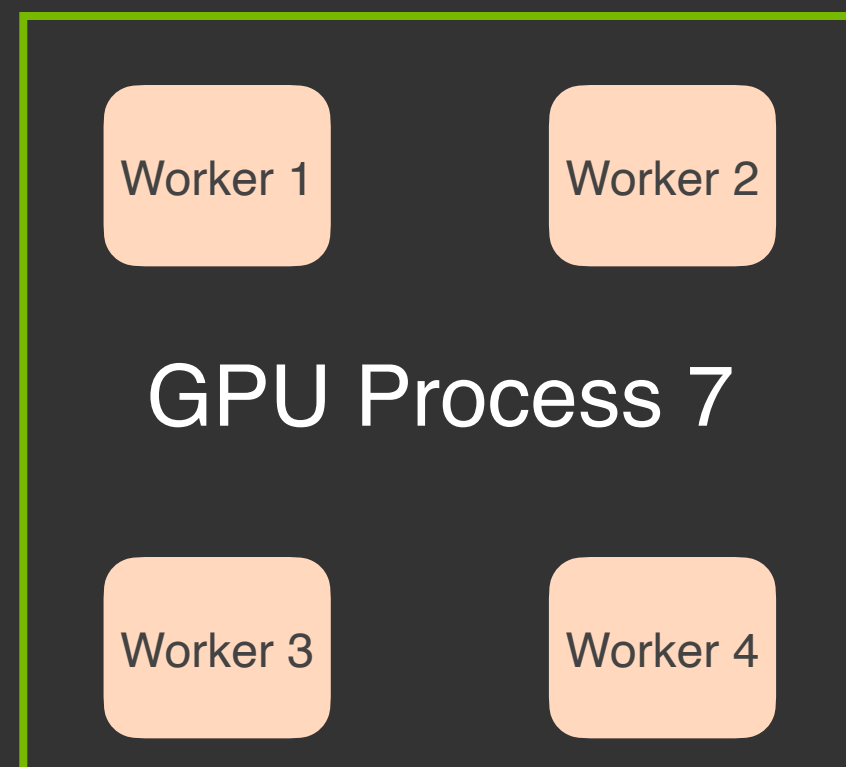
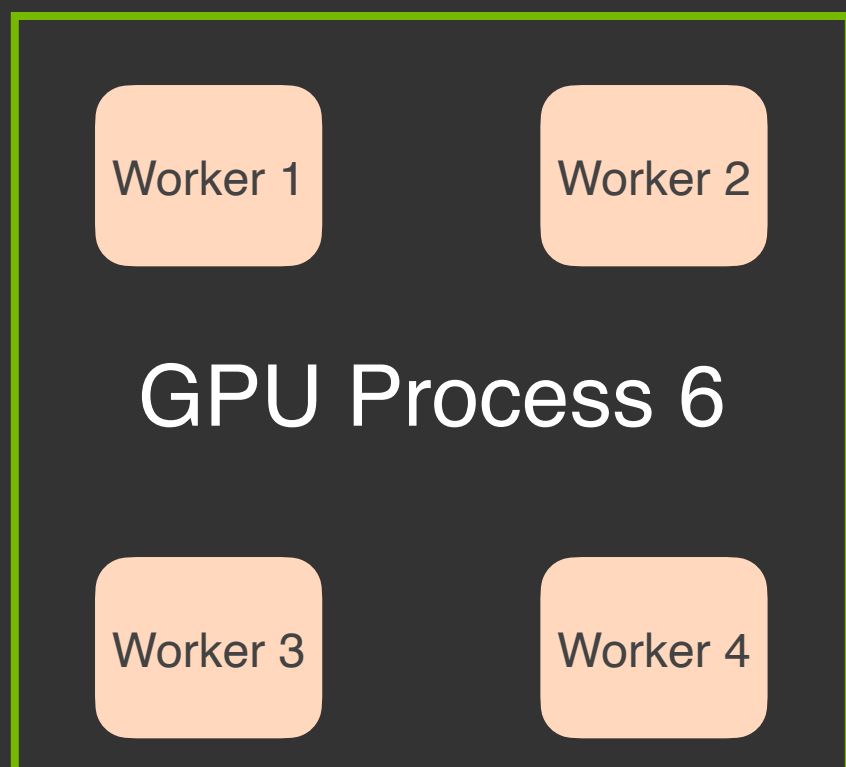
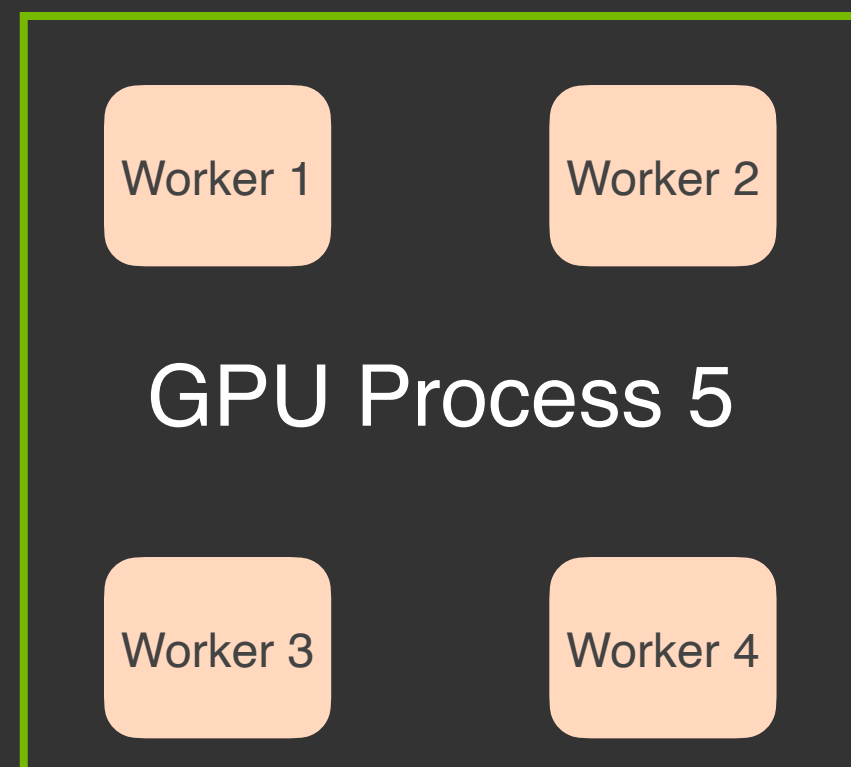
GPU Process 6

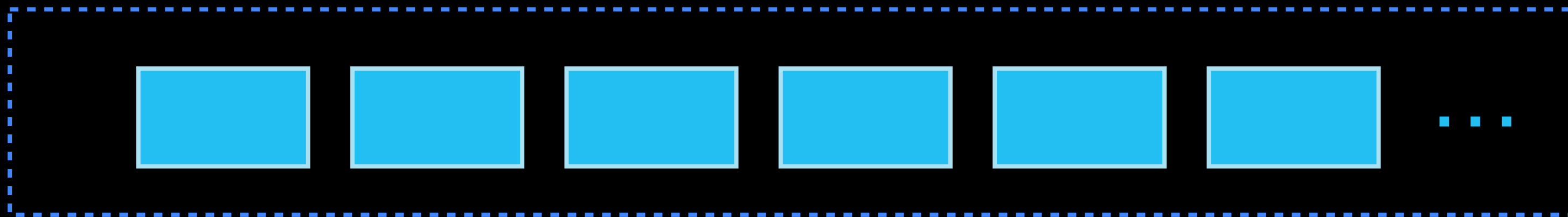
GPU Process 7

GPU Process 8

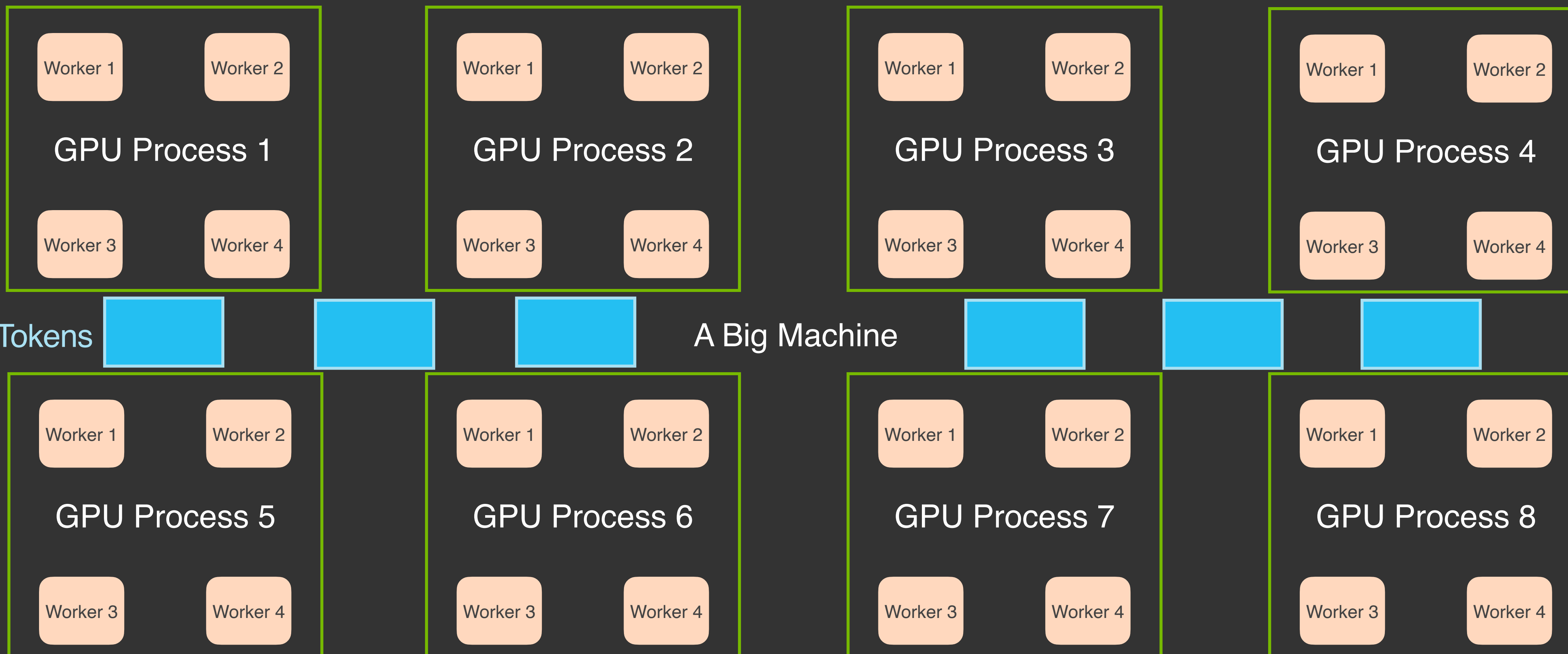


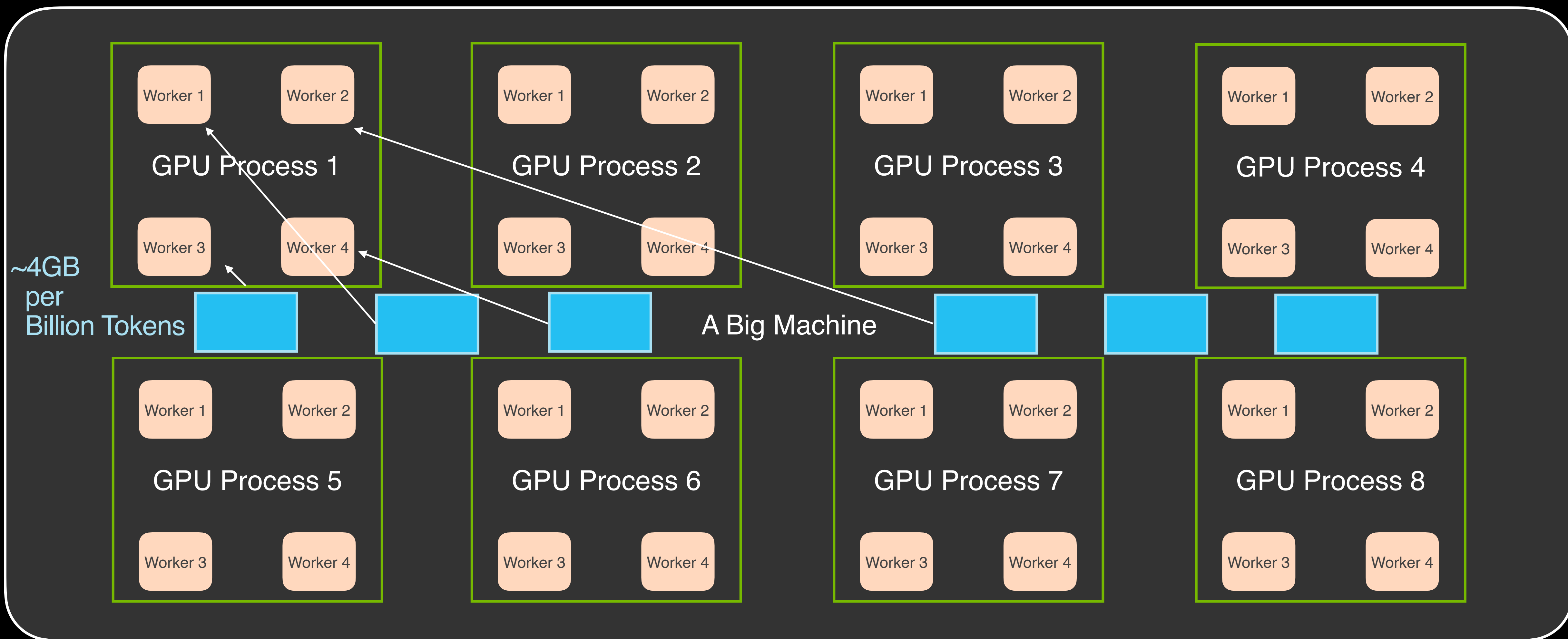
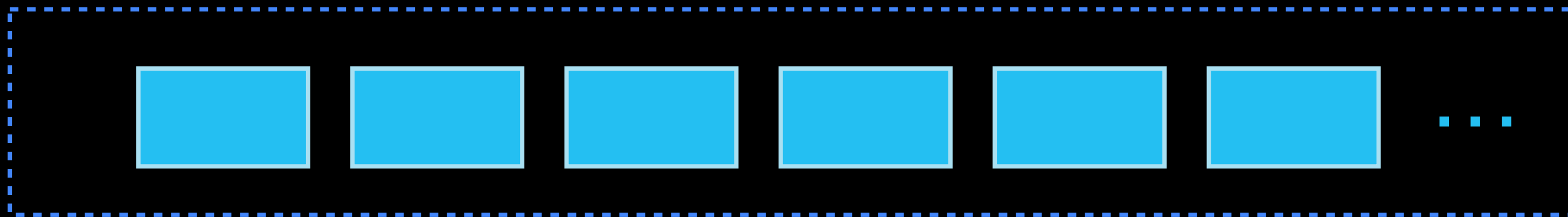
A Big Machine



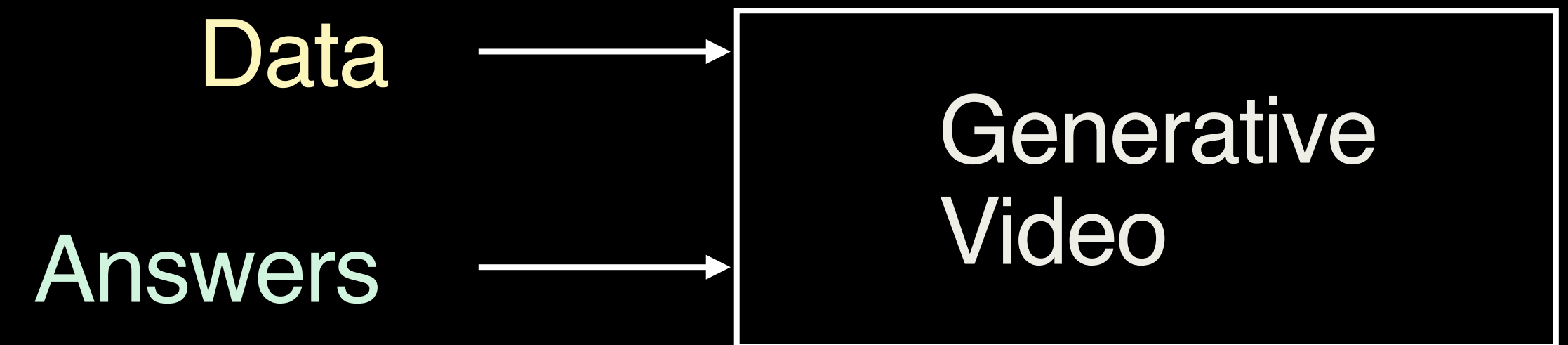


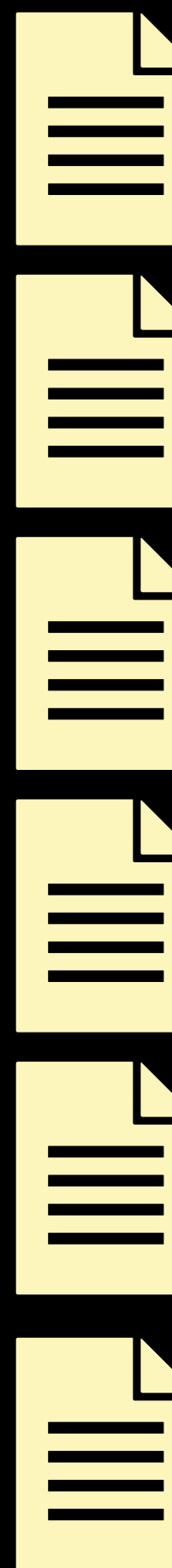
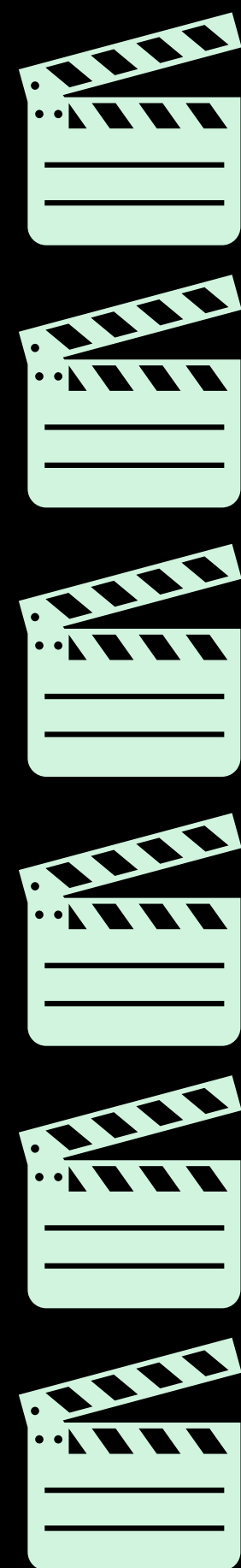
~4GB
per
Billion Tokens





runway



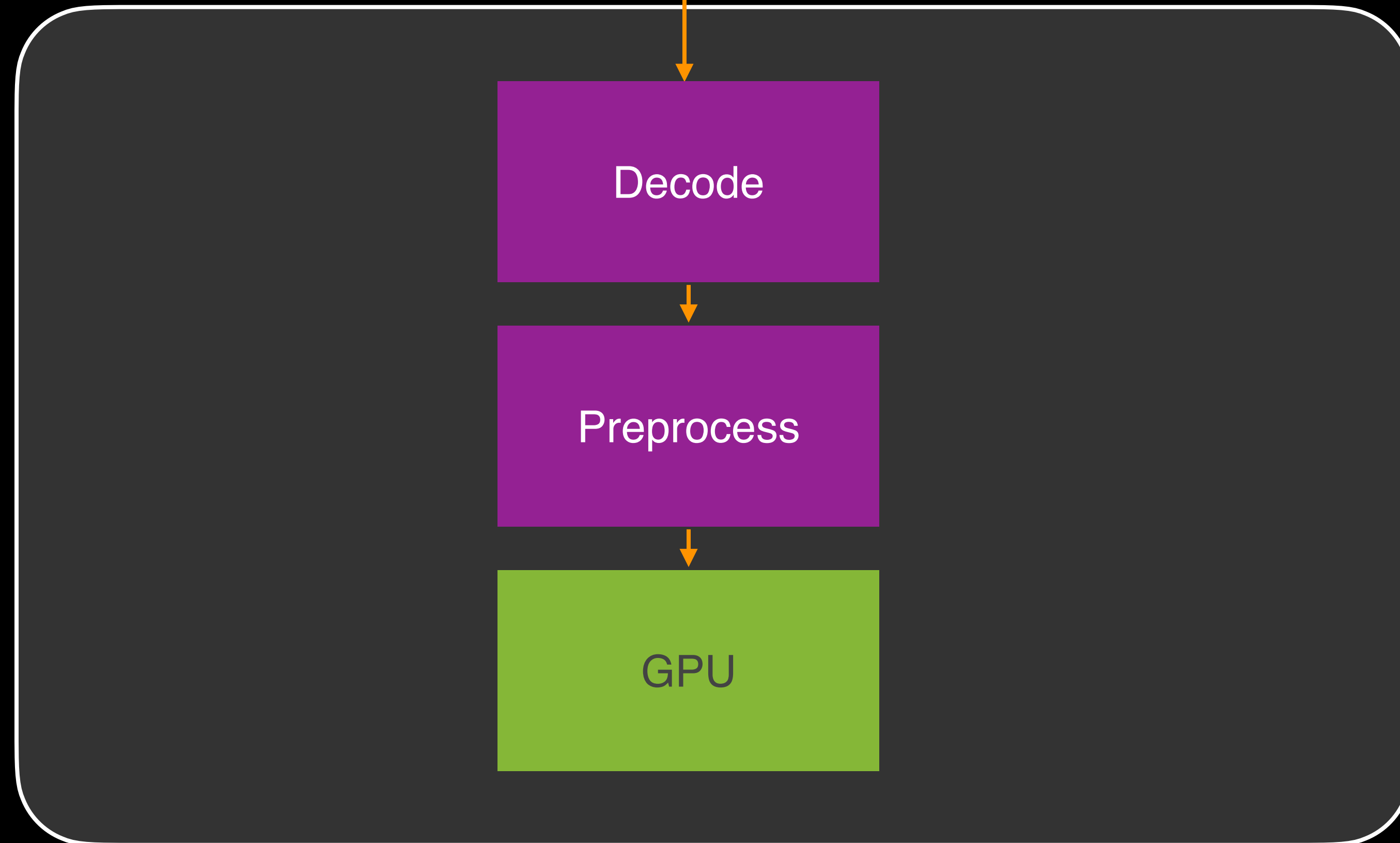


Generative Video

Videos >> Text

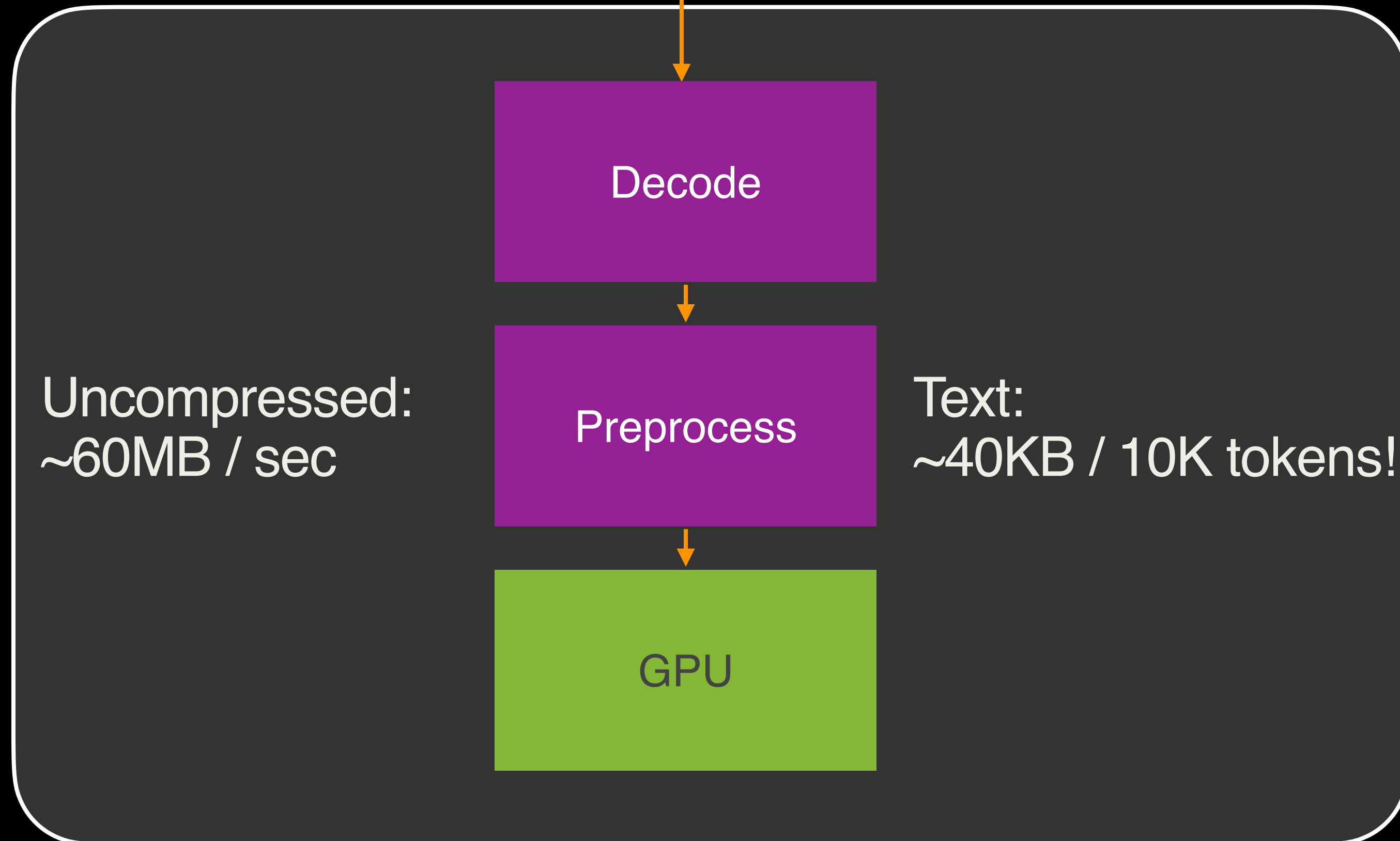


Compressed data size:
 $O(10\text{KB})$ / 1080p frame
 $O(1\text{MB})$ / sec





Compressed data size:
 $O(10\text{KB})$ / 1080p frame
 $O(1\text{MB})$ / sec





Compressed data size:
 $O(10\text{KB})$ / 1080p frame
 $O(1\text{MB})$ / sec

Decode

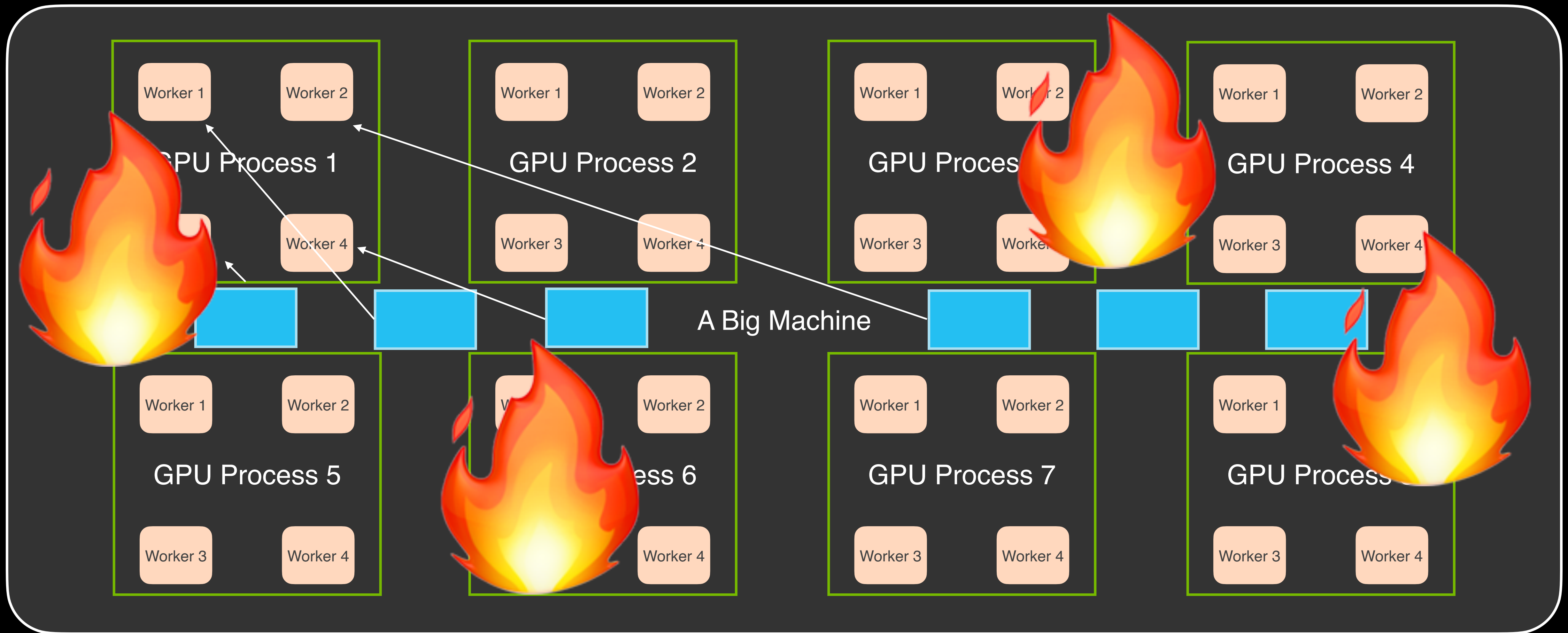
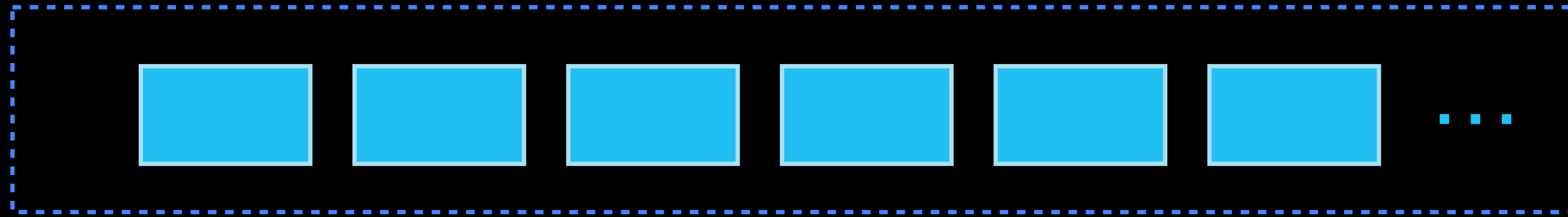
Preprocess

GPU

Uncompressed:
 $\sim 60\text{MB}$ / sec

Text:
 $\sim 40\text{KB}$ / 10K tokens!

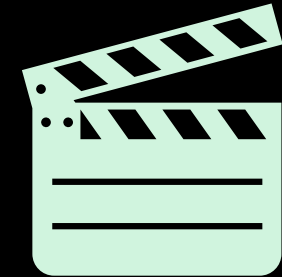
Don't Bottleneck the GPU!



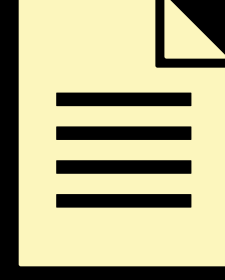
Data Refining



Videos



Prompts



Data

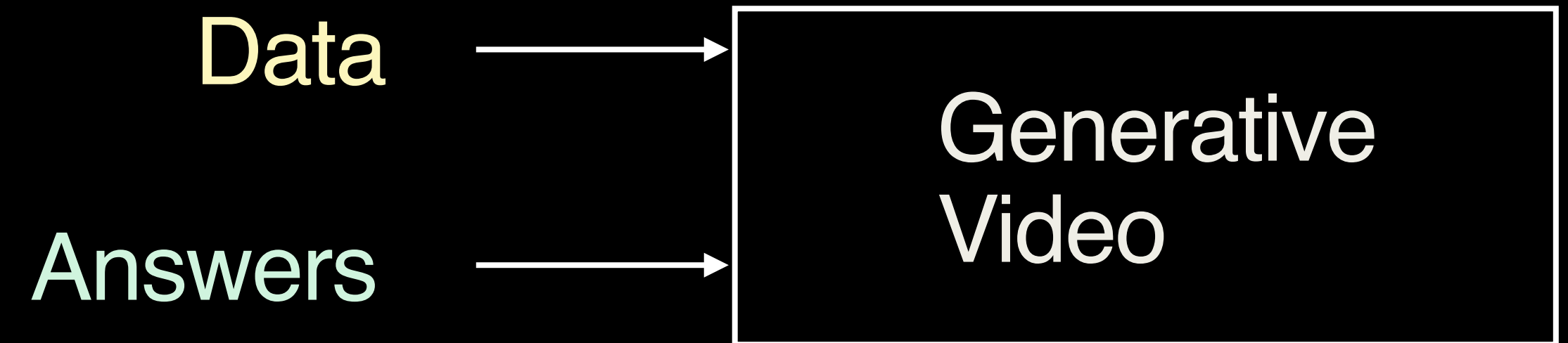
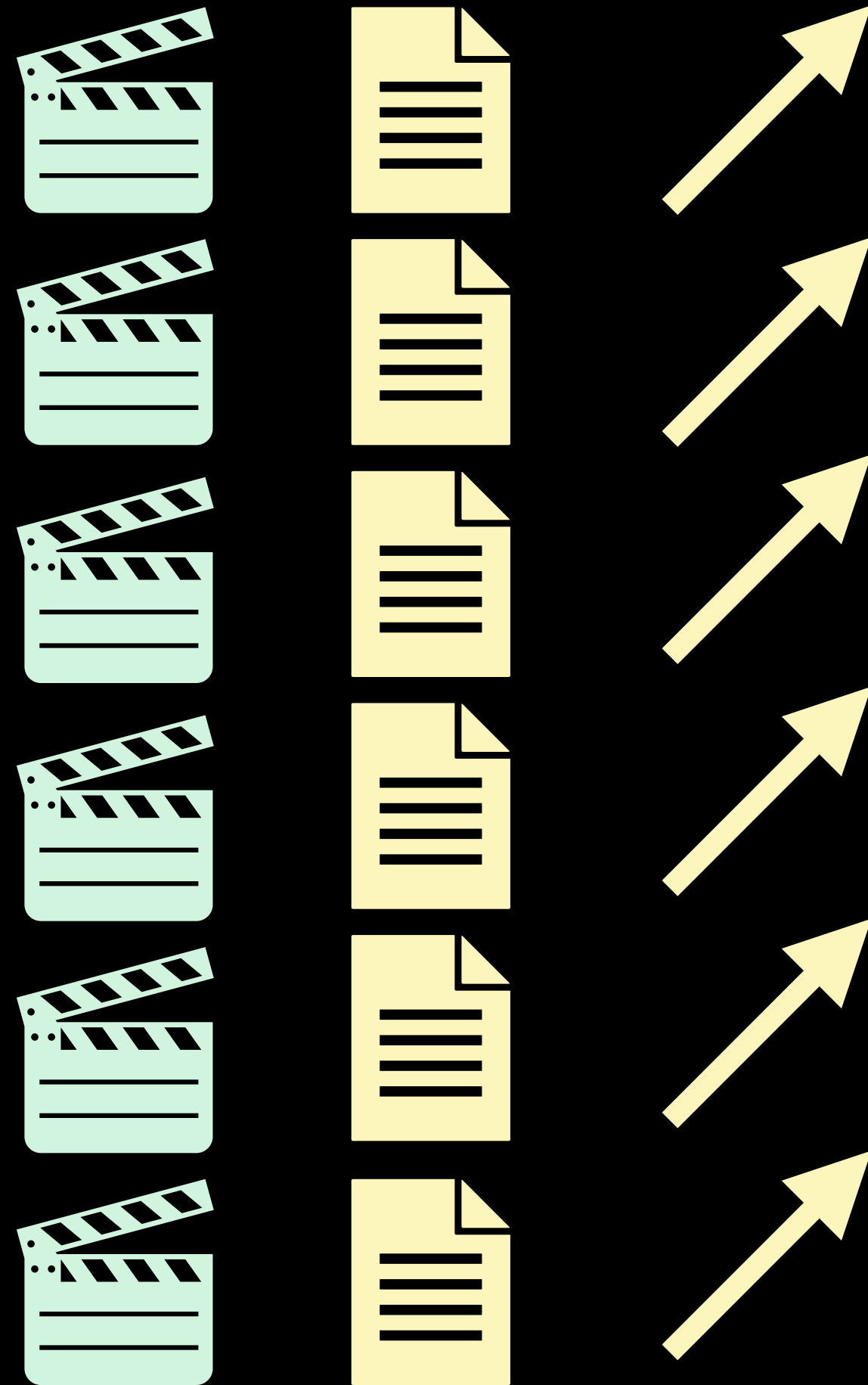


Answers



Generative
Video

Videos Prompts Embeddings



Structure and Content-Guided Video Synthesis with Diffusion Models

Patrick Esser Johnathan Chiu Parmida Atighehchian
Jonathan Granskog Anastasis Germanidis
Runway

<https://research.runwayml.com/gen1>

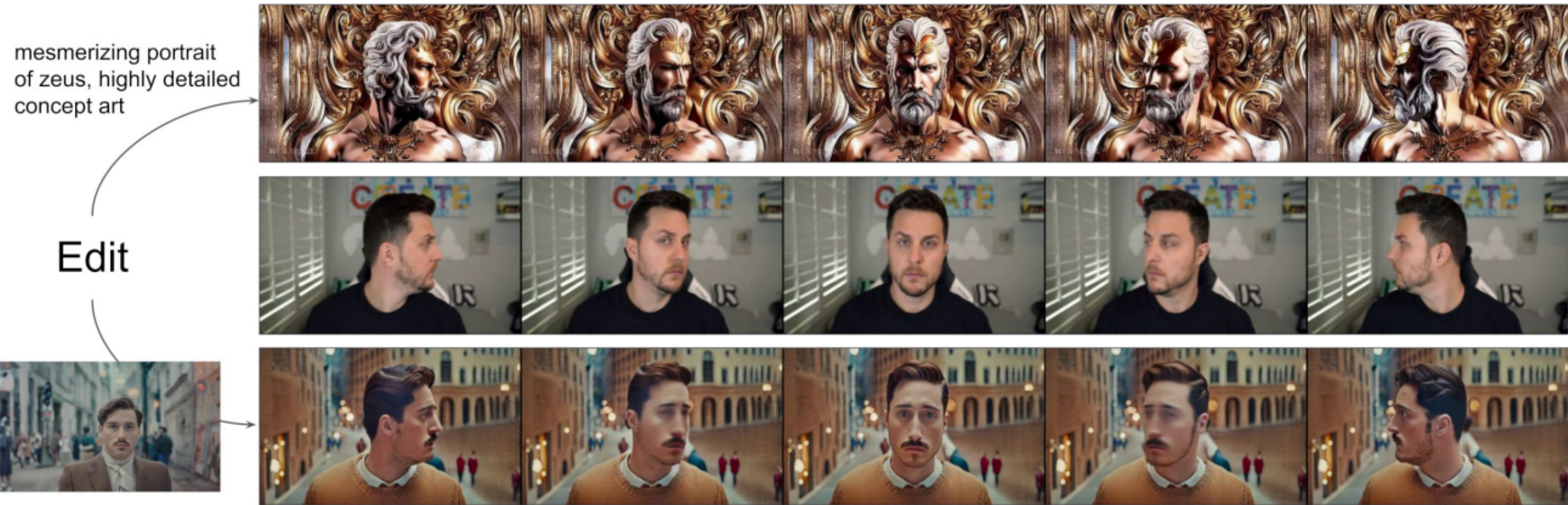


Figure 1. **Guided Video Synthesis** We present an approach based on latent video diffusion models that synthesizes videos (top and bottom) guided by content described through text (top) or images (bottom) while keeping the structure of an input video (middle).

Structure and Content-Guided Video Synthesis with Diffusion Models

Patrick Esser Johnathan Chiu Parmida Atighehchian
Jonathan Granskog Anastasis Germanidis
Runway

<https://research.runwayml.com/gen1>

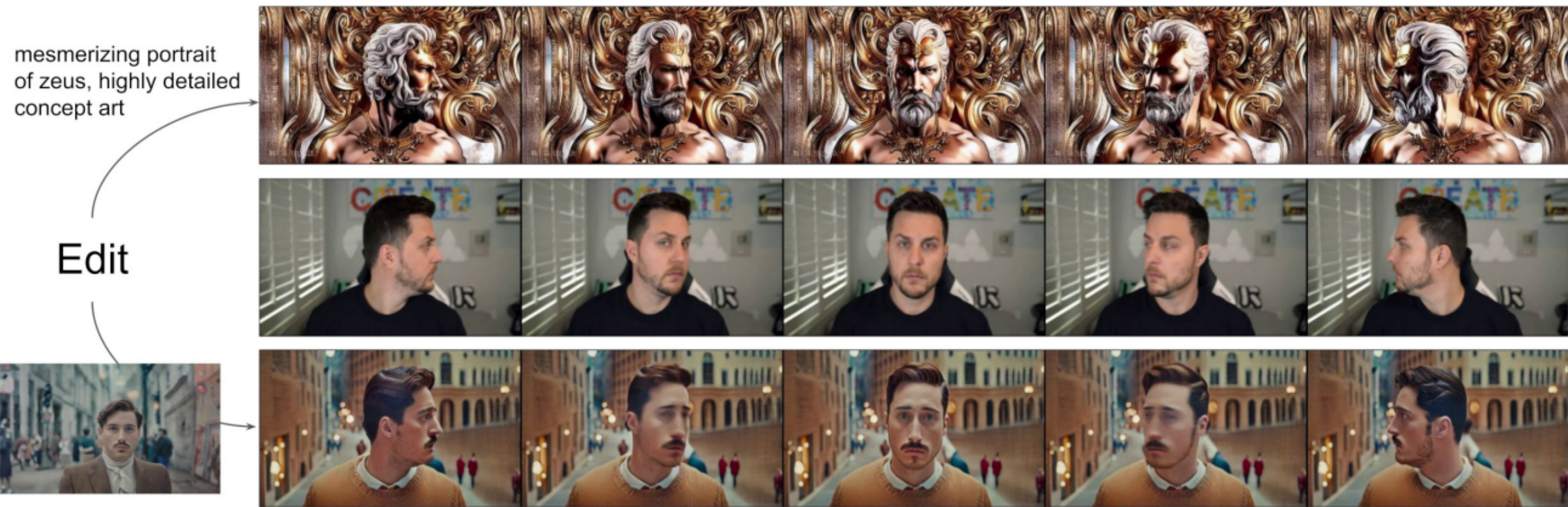
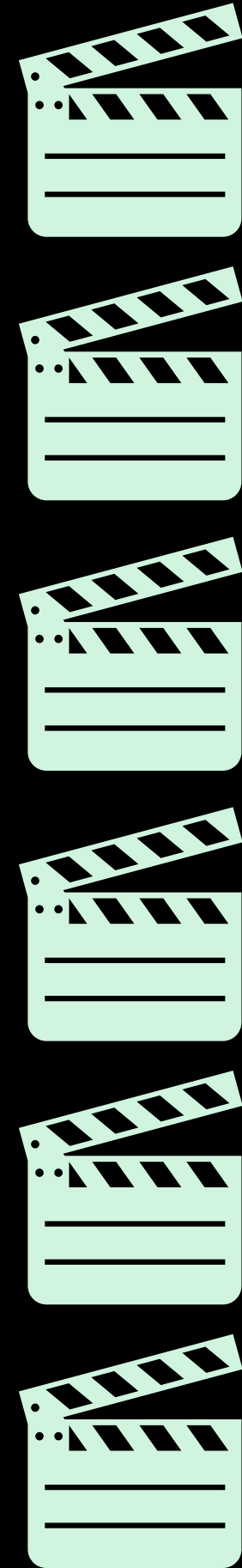


Figure 1. **Guided Video Synthesis** We present an approach based on latent video diffusion models that synthesizes videos (top and bottom) guided by content described through text (top) or images (bottom) while keeping the structure of an input video (middle).

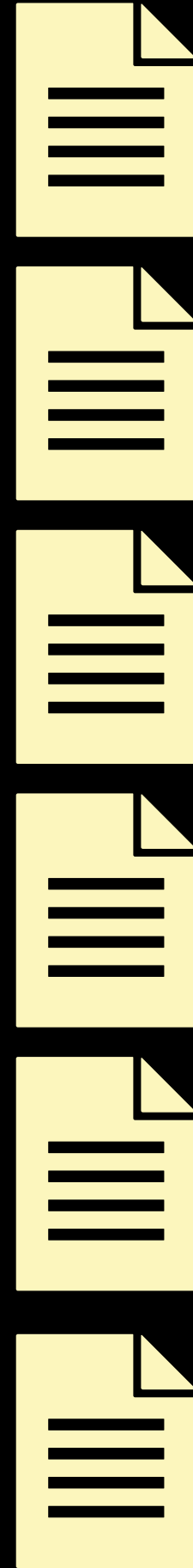


Depth Anything V2, Neurips '24

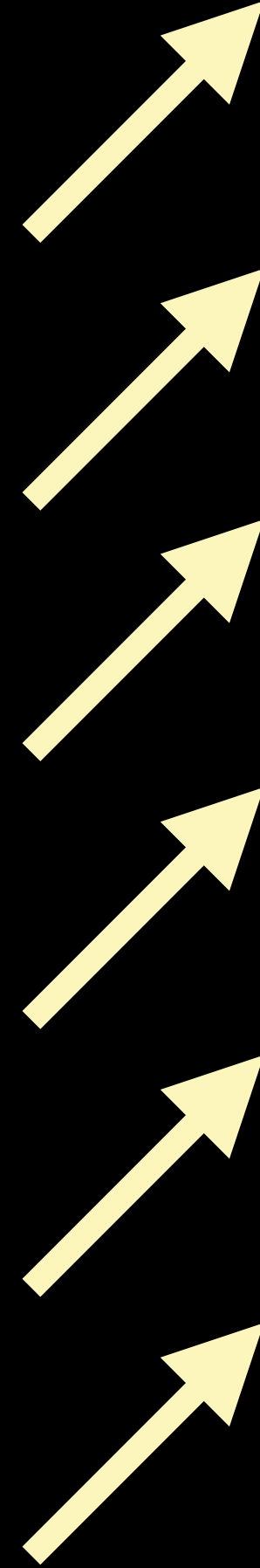
Videos



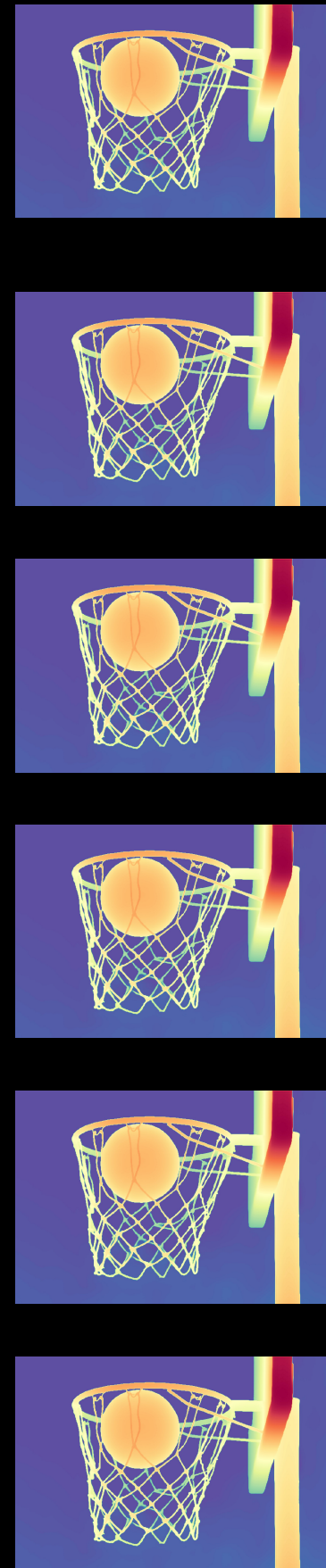
Prompts



Embeddings



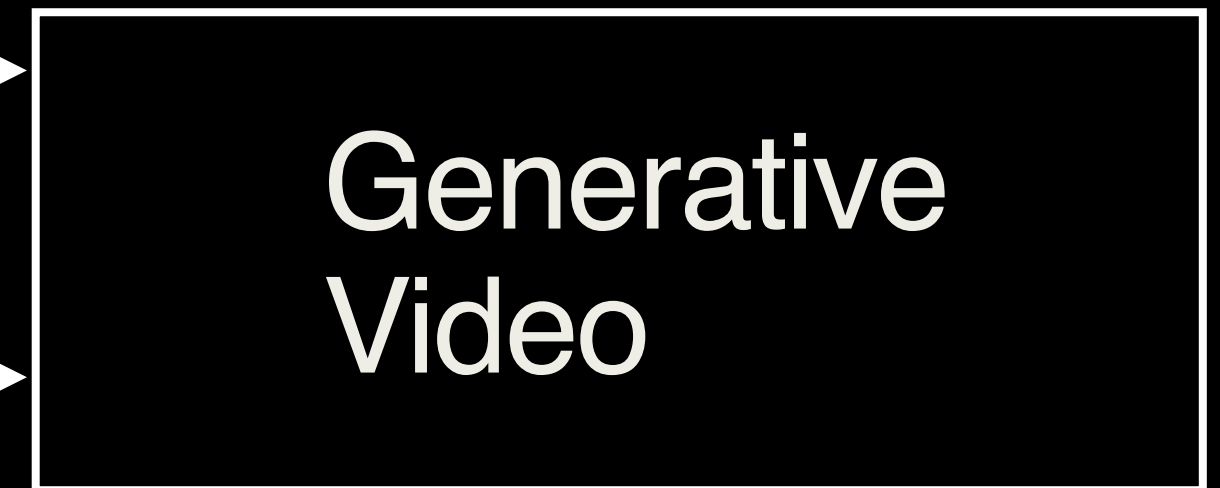
Depth Maps



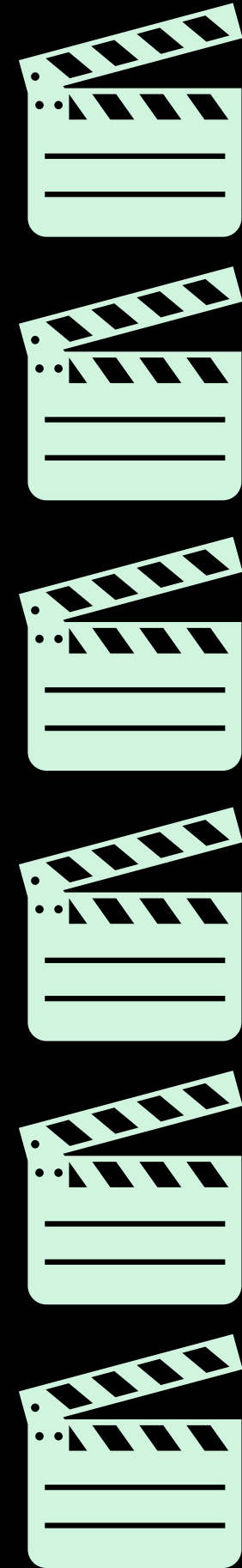
Data
Answers



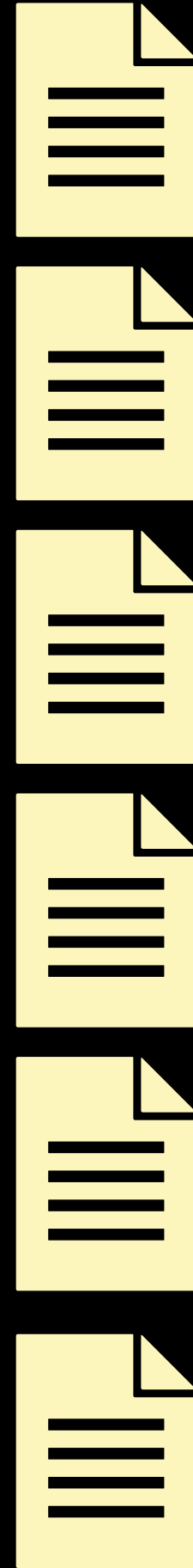
Generative
Video



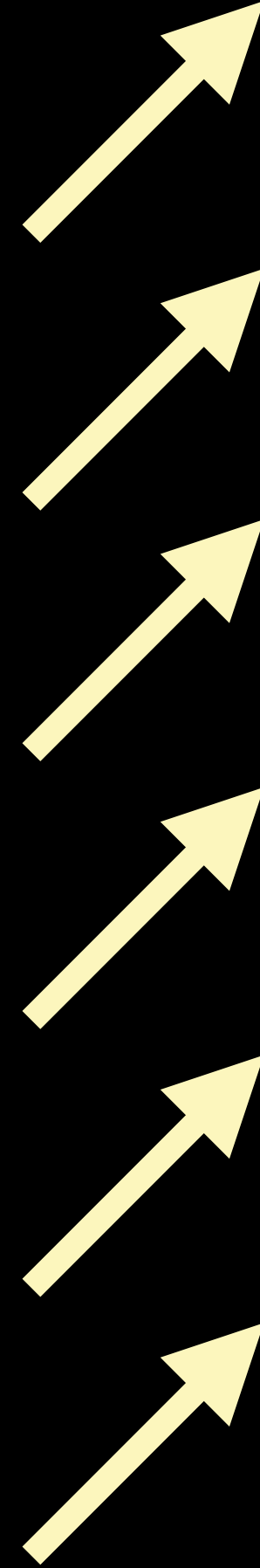
Videos



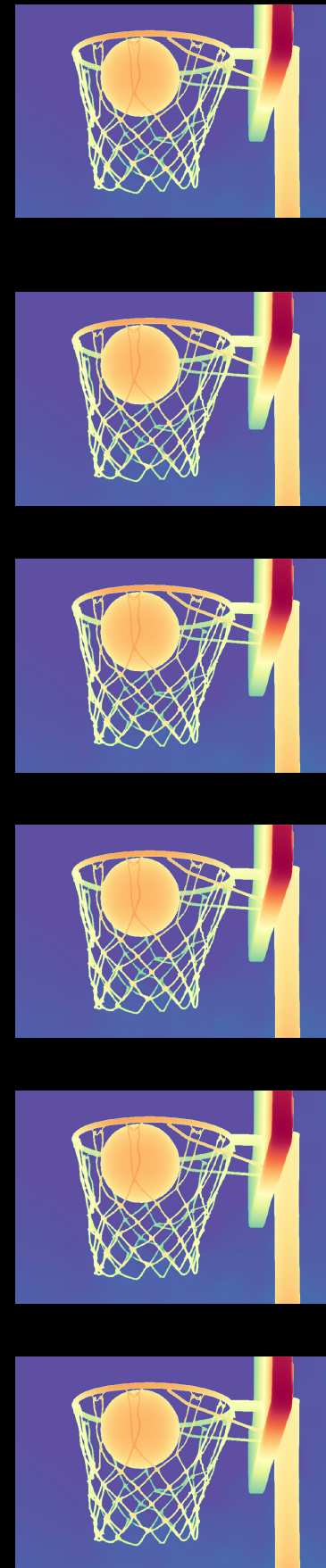
Prompts



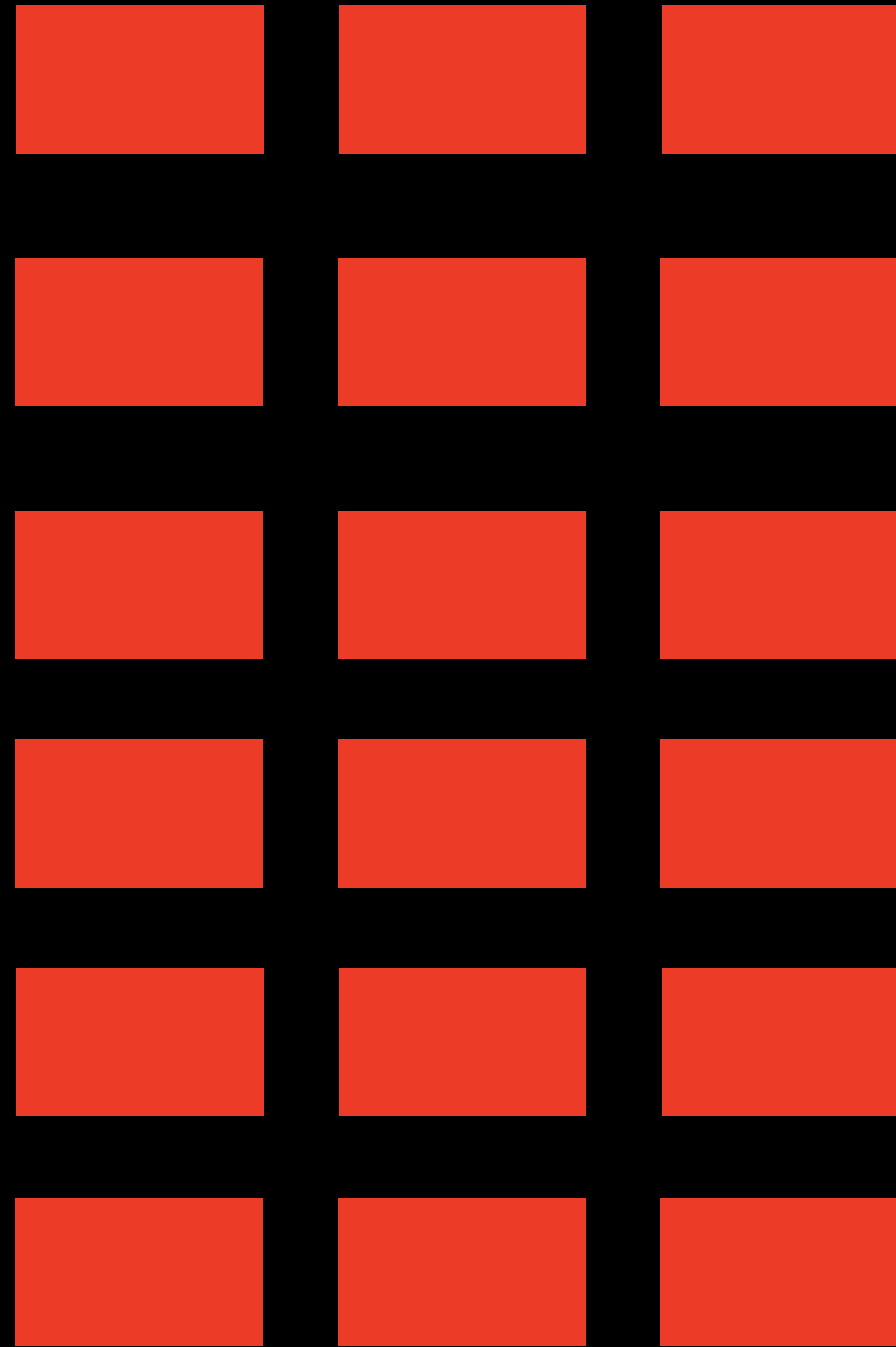
Embeddings



Depth Maps

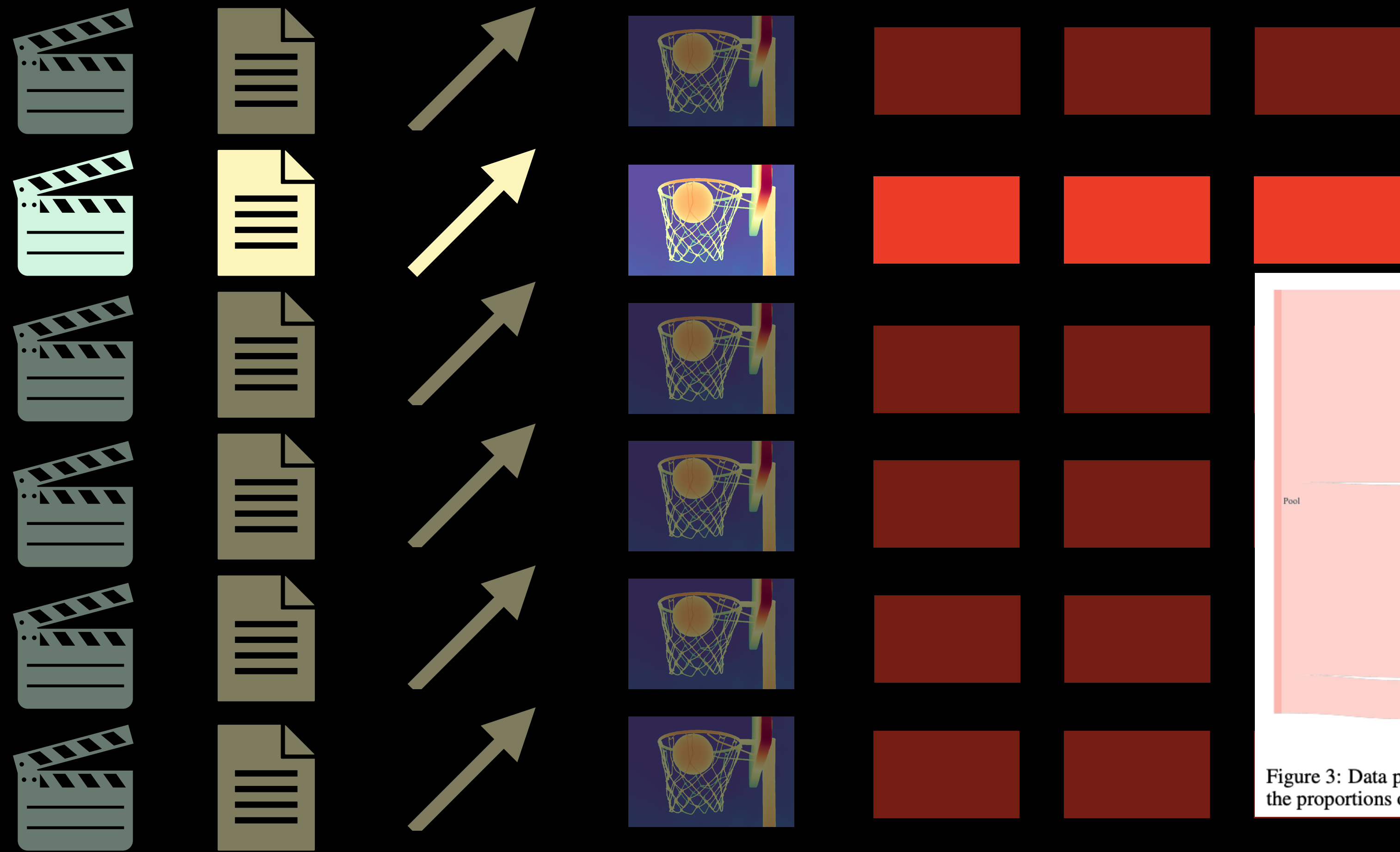


resolution, fps, duration, ...



Generative
Video

Videos Prompts Embeddings Depth Maps resolution, fps, duration, ...



Wan 2.1, Alibaba '25

Taking Stock

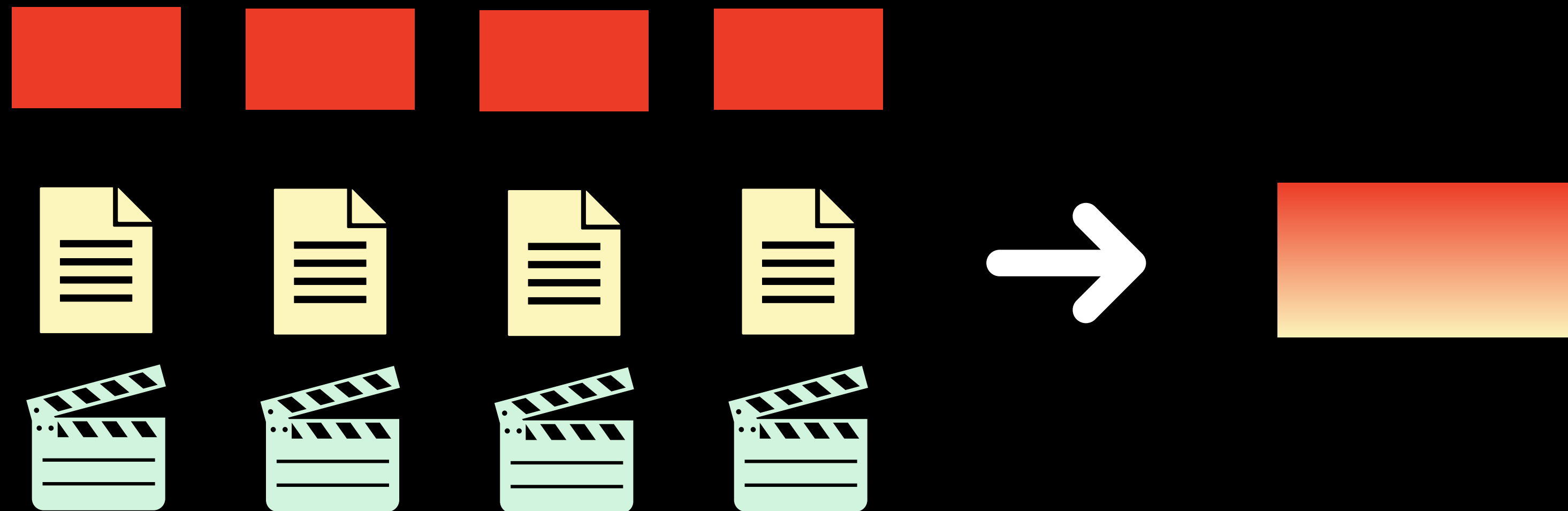
1. Massive unstructured data
2. Structured data
3. Large column appends w/ backfills
4. Dynamic partitioning across workers

One Off vs. Continual Iteration

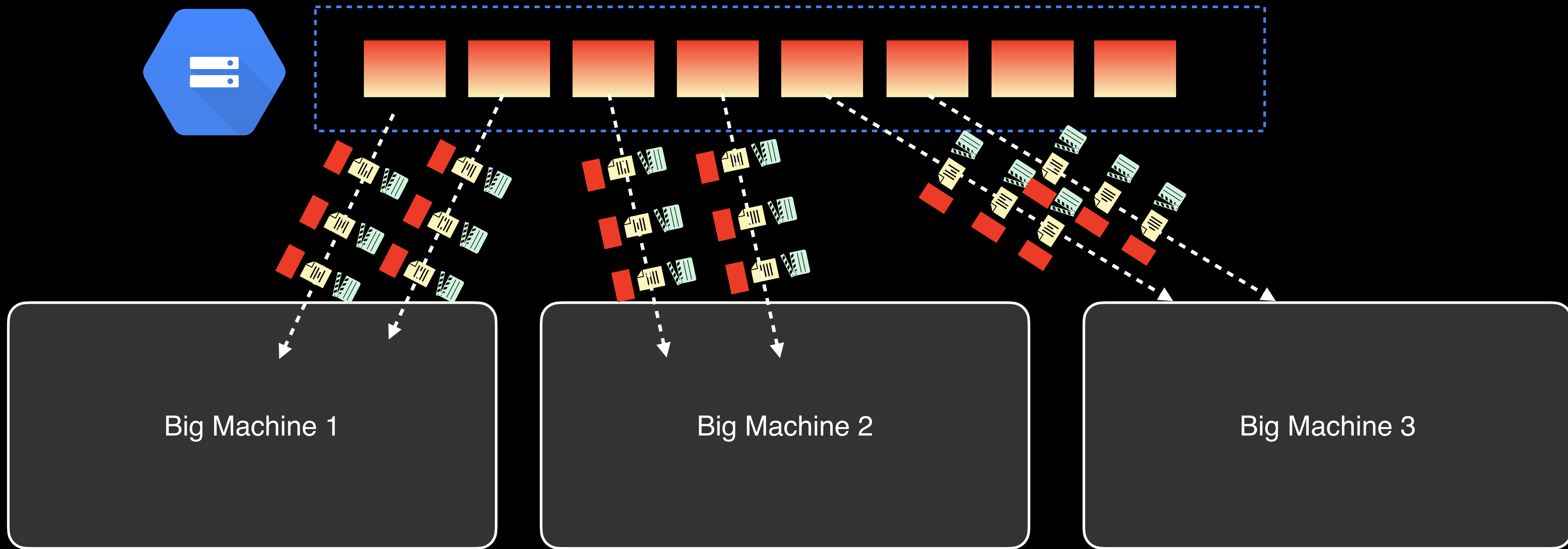
Taking Stock

1. Massive unstructured data
2. Structured data
3. Large column appends w/ backfills
4. Dynamic partitioning across workers
5. EDA

Webdataset (aka tarballs + API)



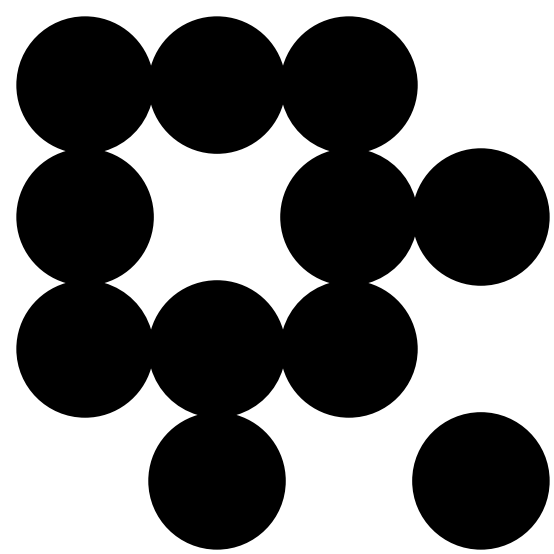
Webdataset (aka tarballs + API)



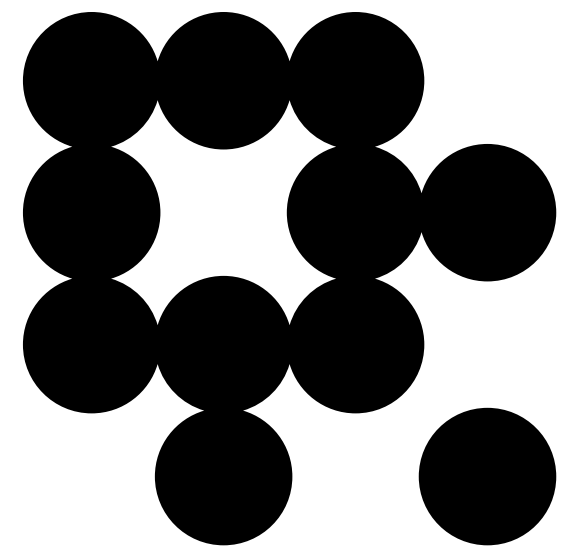
Taking Stock — Webdataset

- ✓ 1. Massive unstructured data
- ✓ 2. Structured data
- ✗ 3. Large column appends w/ backfills
- 😐 4. Dynamic partitioning across workers
- ✗ 5. EDA

Object storage-native
Columnar
Heterogeneous types
Fast random access

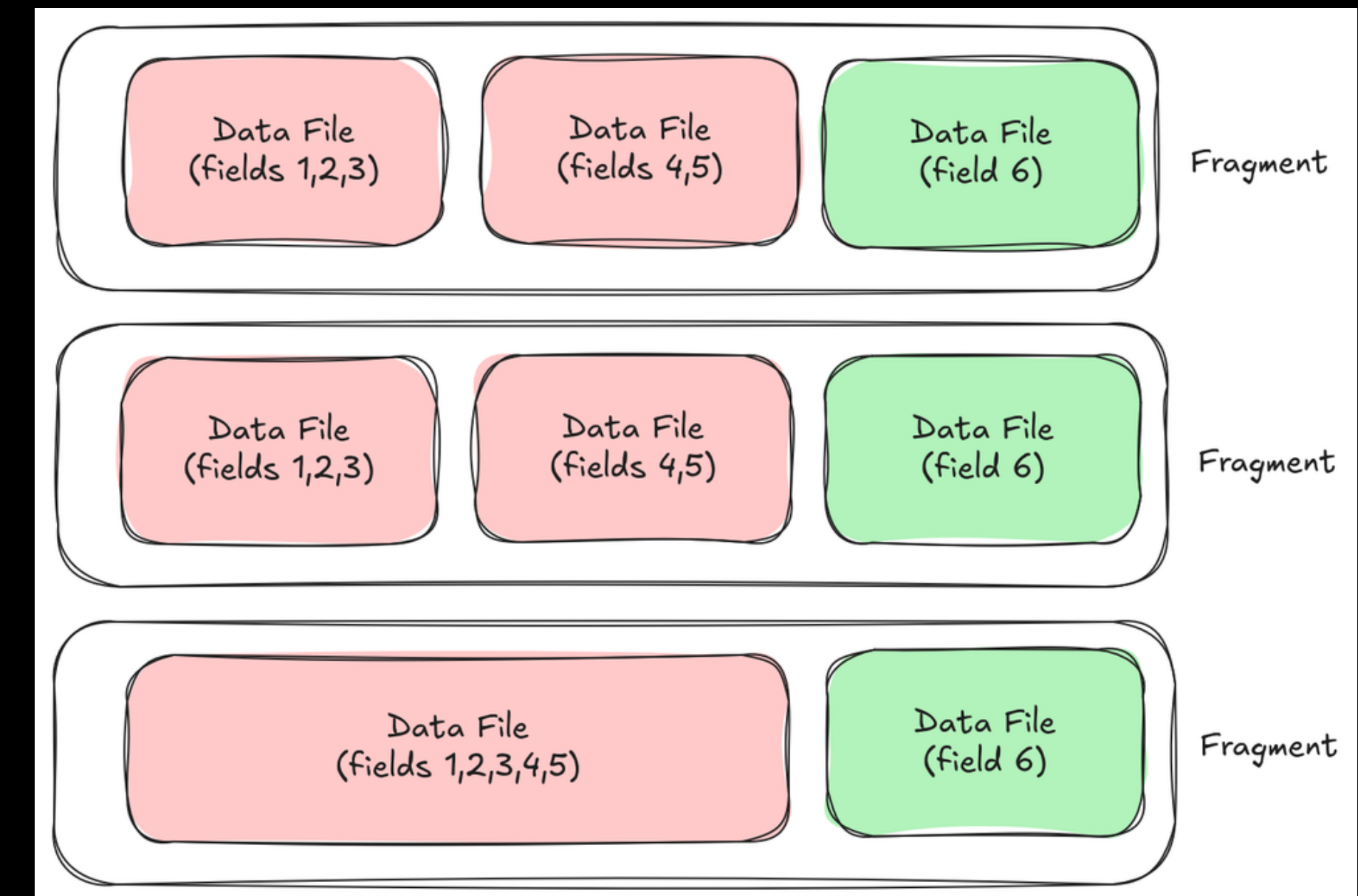


LanceDB



LanceDB

- File + table format
- Column appends w/o rewriting data
- Fast random access
- Multimodal support
- Lots of other bells and whistles (versioning, arrow integration, vector search, etc...)



Taking Stock — Lance

- ✓ 1. Massive unstructured data
- ✓ 2. Structured data
- ✓ 3. Column appends
- 4. Dynamic partitioning across workers
- 5. EDA

Lance Dynamic Partitioning



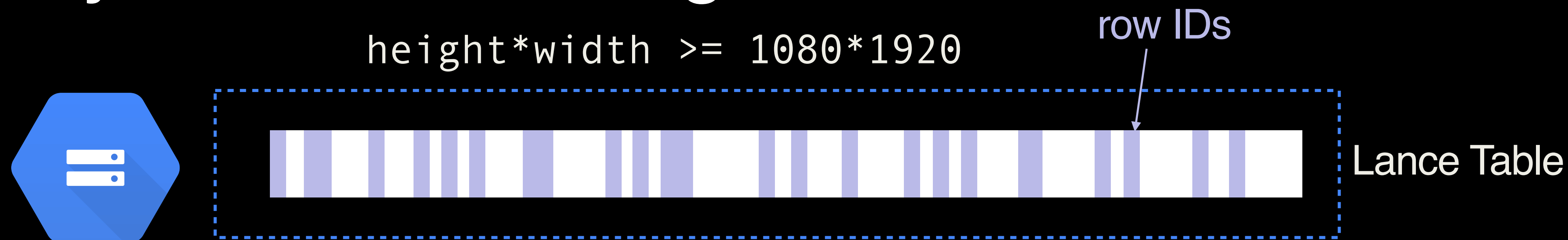
Lance Table

Big Machine 1

Big Machine 2

Big Machine 3

Lance Dynamic Partitioning

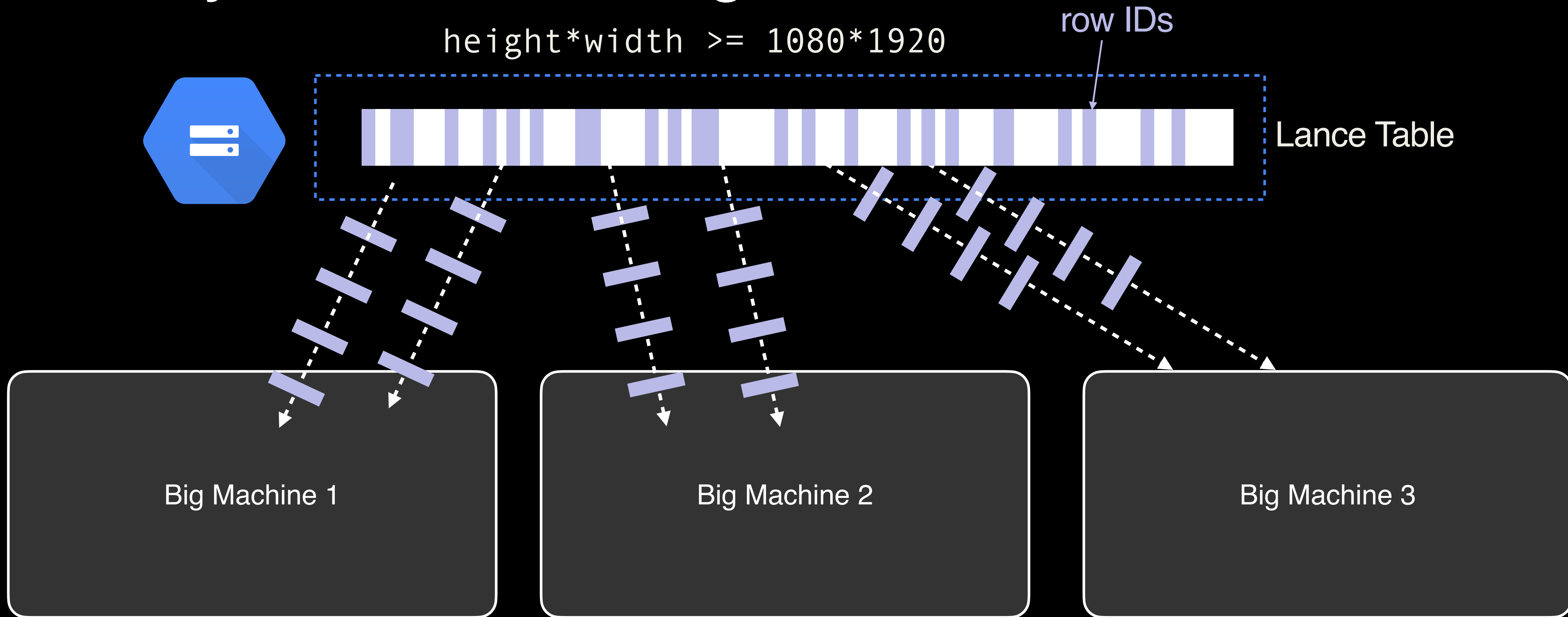


Big Machine 1

Big Machine 2

Big Machine 3

Lance Dynamic Partitioning



Taking Stock — Lance

- ✓ 1. Massive unstructured data
- ✓ 2. Structured data
- ✓ 3. Column appends
- ✓ 4. Dynamic partitioning across workers
- 5. EDA

Taking Stock — Lance



- ✓ 1. Massive unstructured data
- ✓ 2. Structured data
- ✓ 3. Column appends
- ✓ 4. Dynamic partitioning across workers
- ✓ 5. EDA

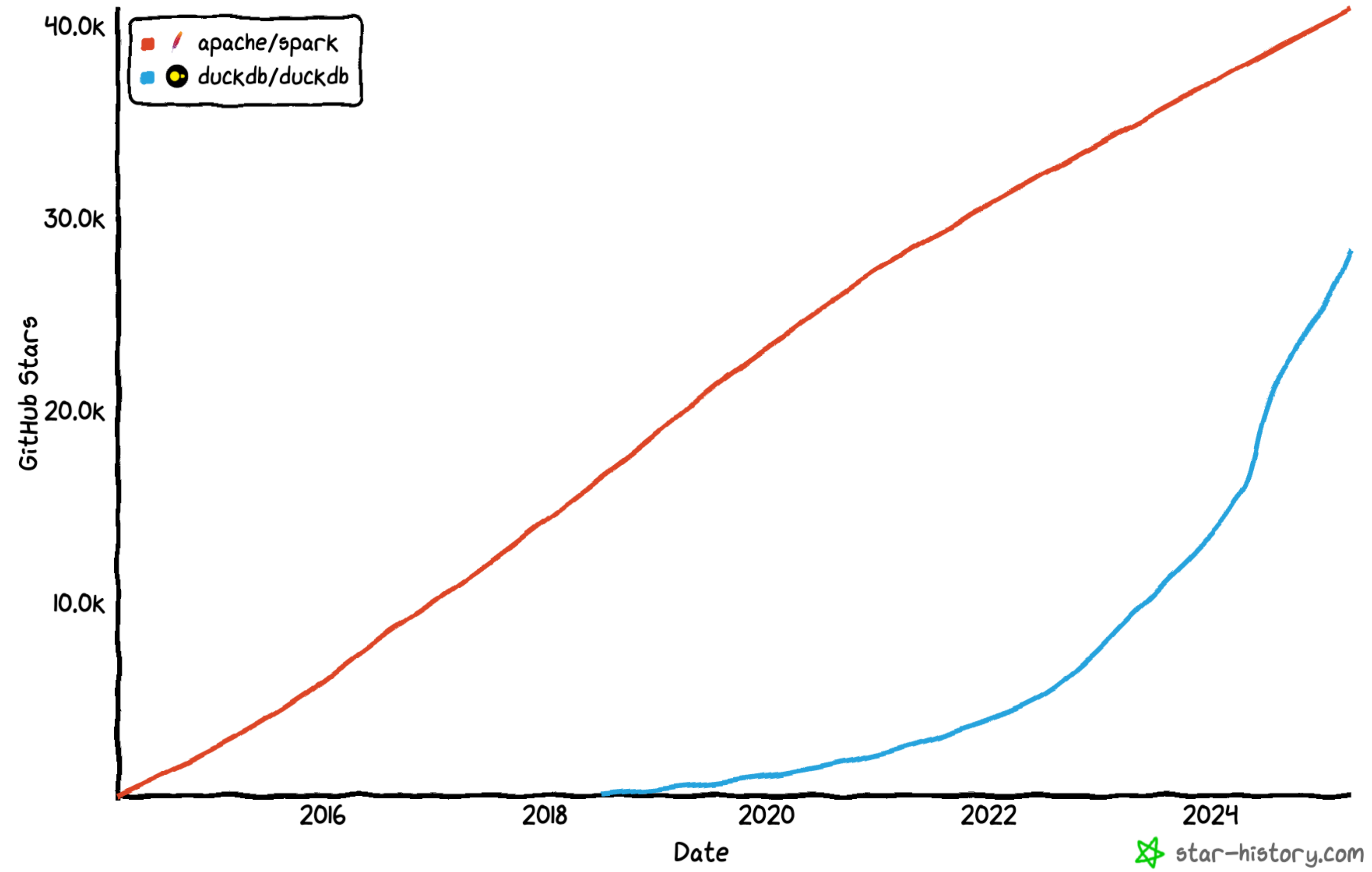
```
SELECT
  COUNT(1) AS num_hd
FROM my_videos
WHERE
  height * width >= 1080 * 1920
```


THE WORLD AFTER YOUR DATA IS IN A DB

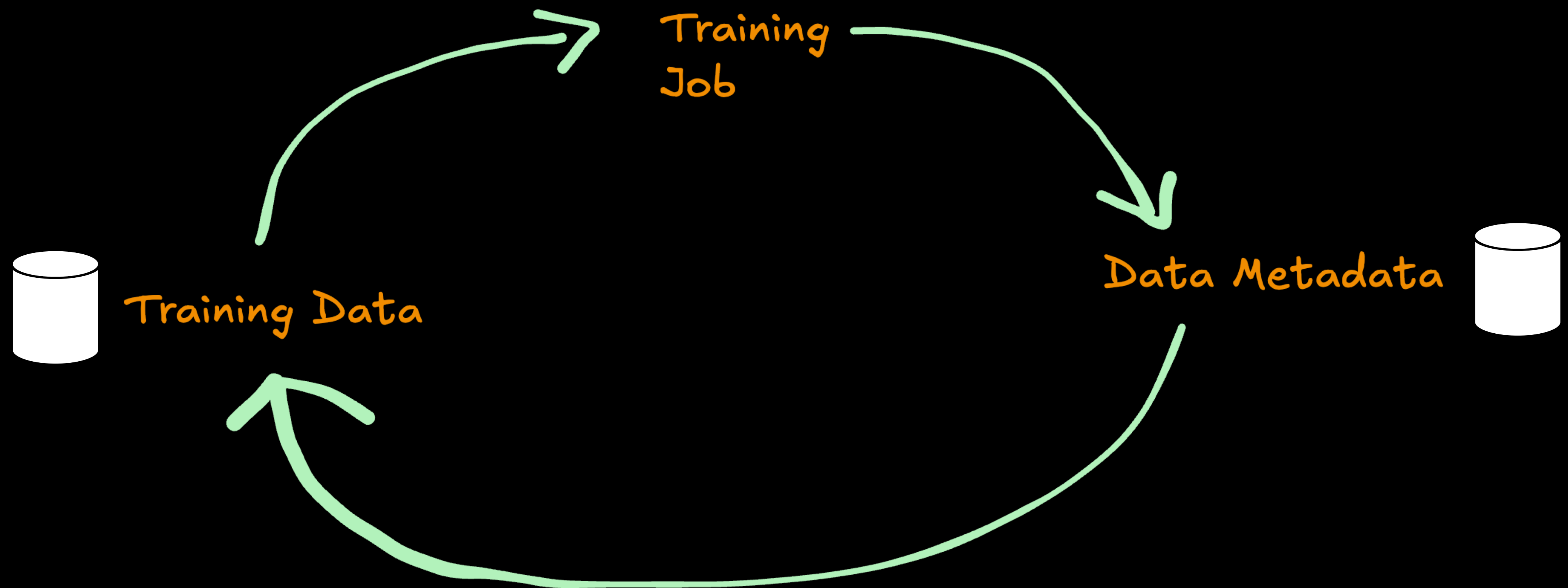


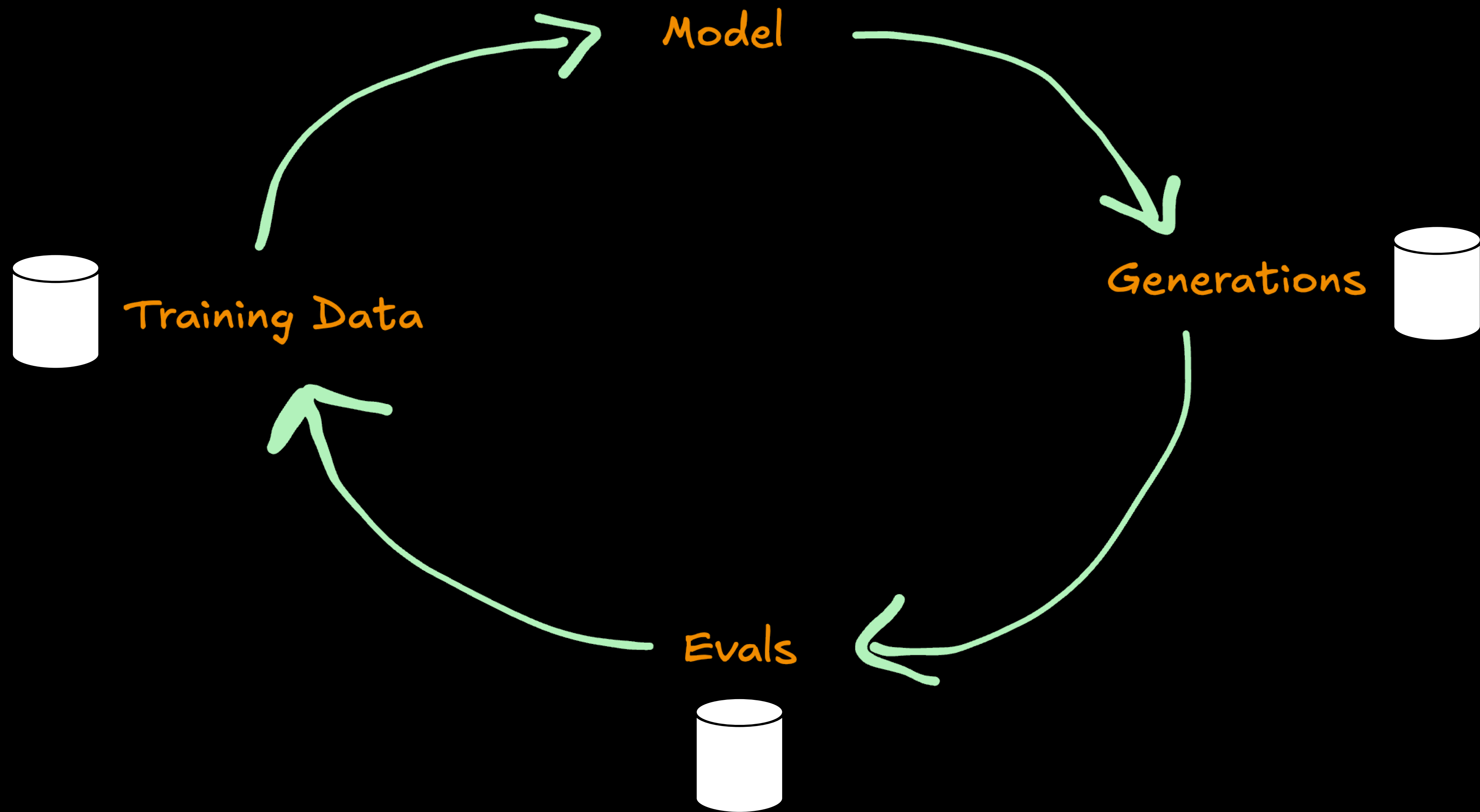
imgflip.com

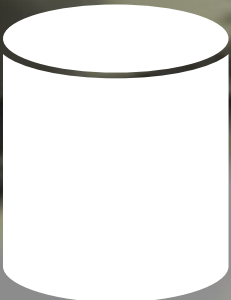
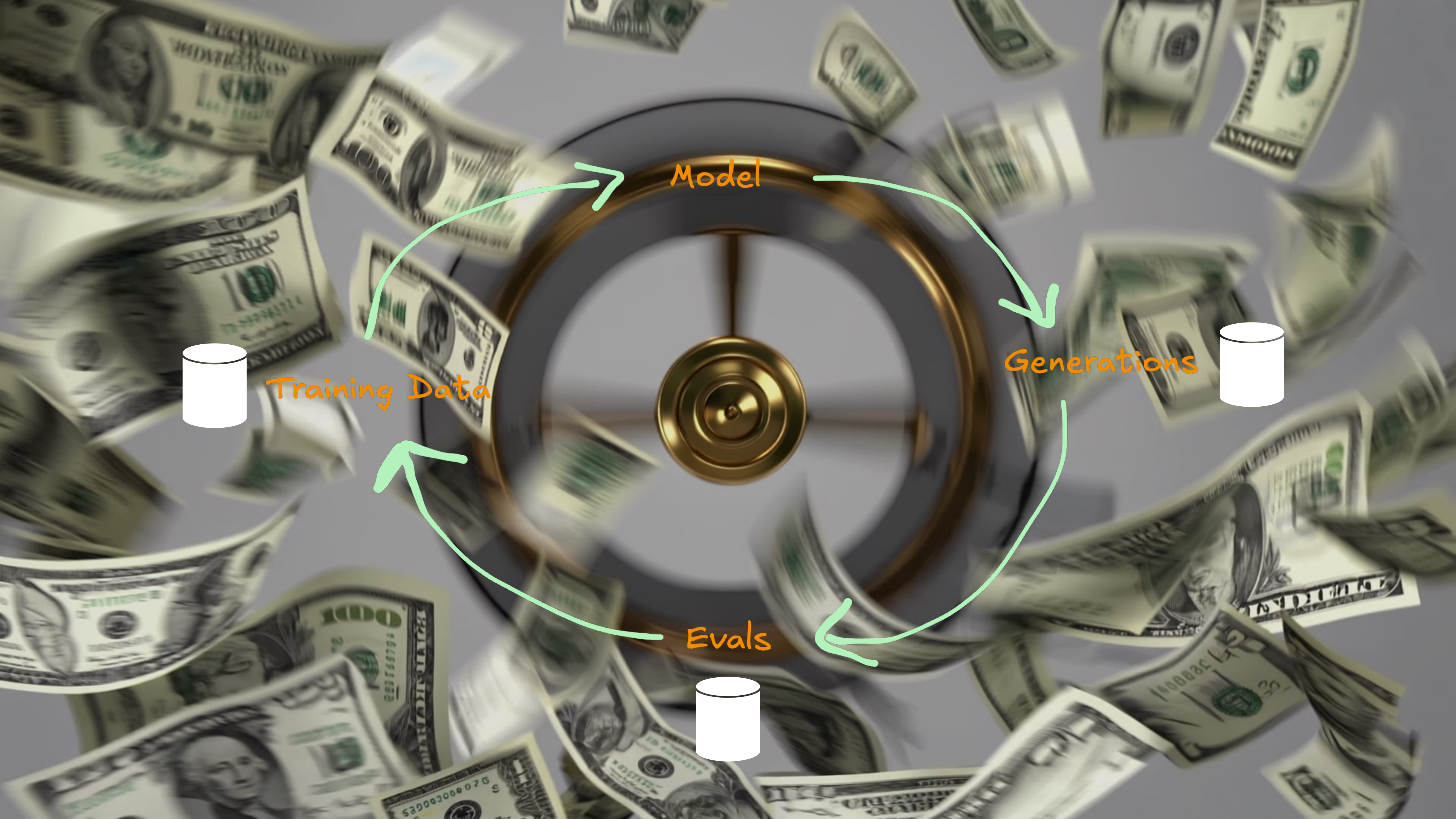
Star History











Training Data

Model



Generations

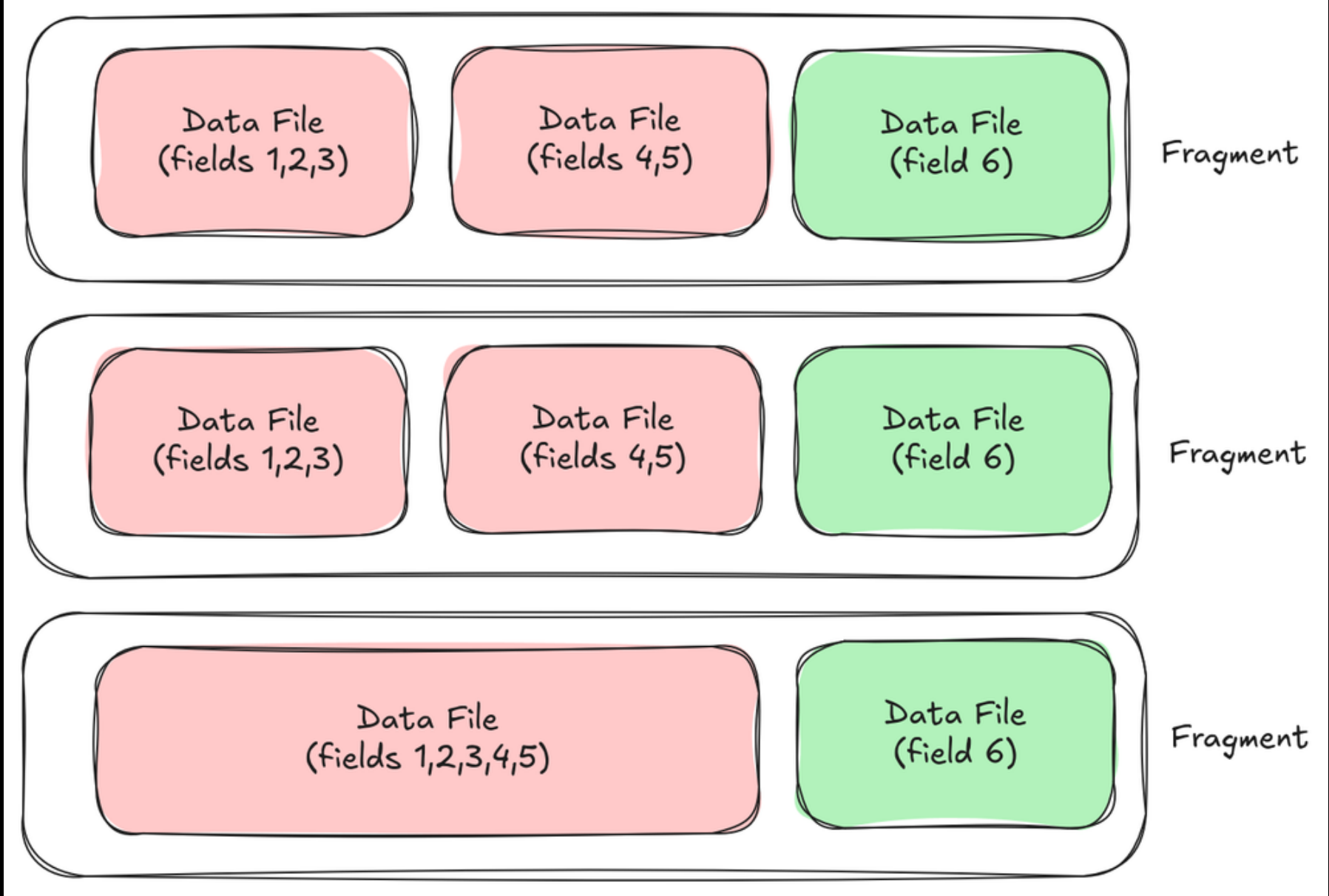


Evals

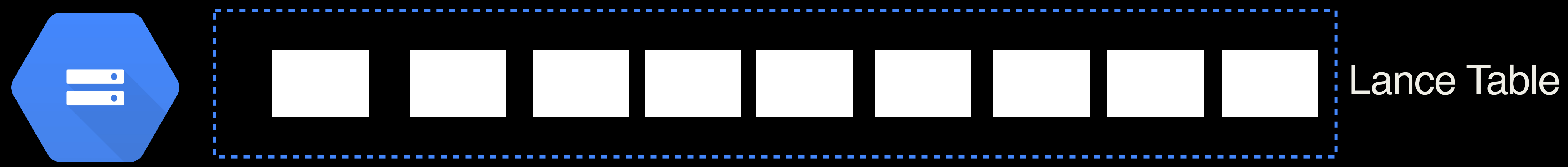
That easy, huh?

Lessons Learned





Lance Fragment Partitioning

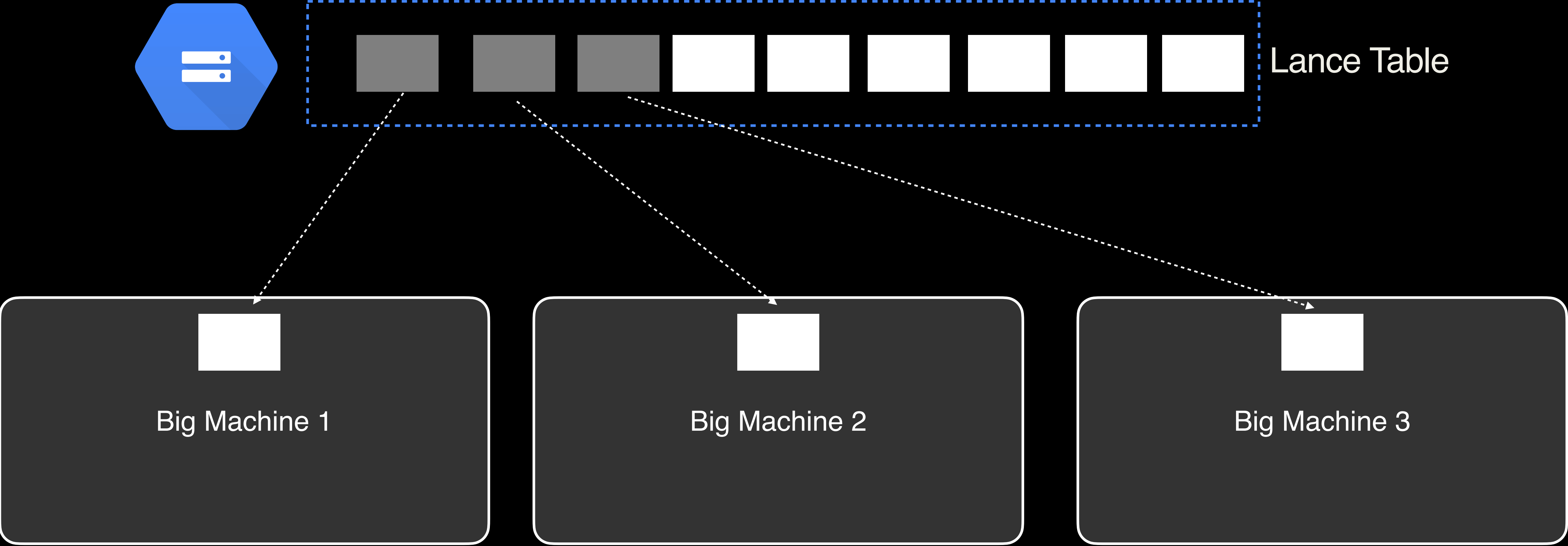


Big Machine 1

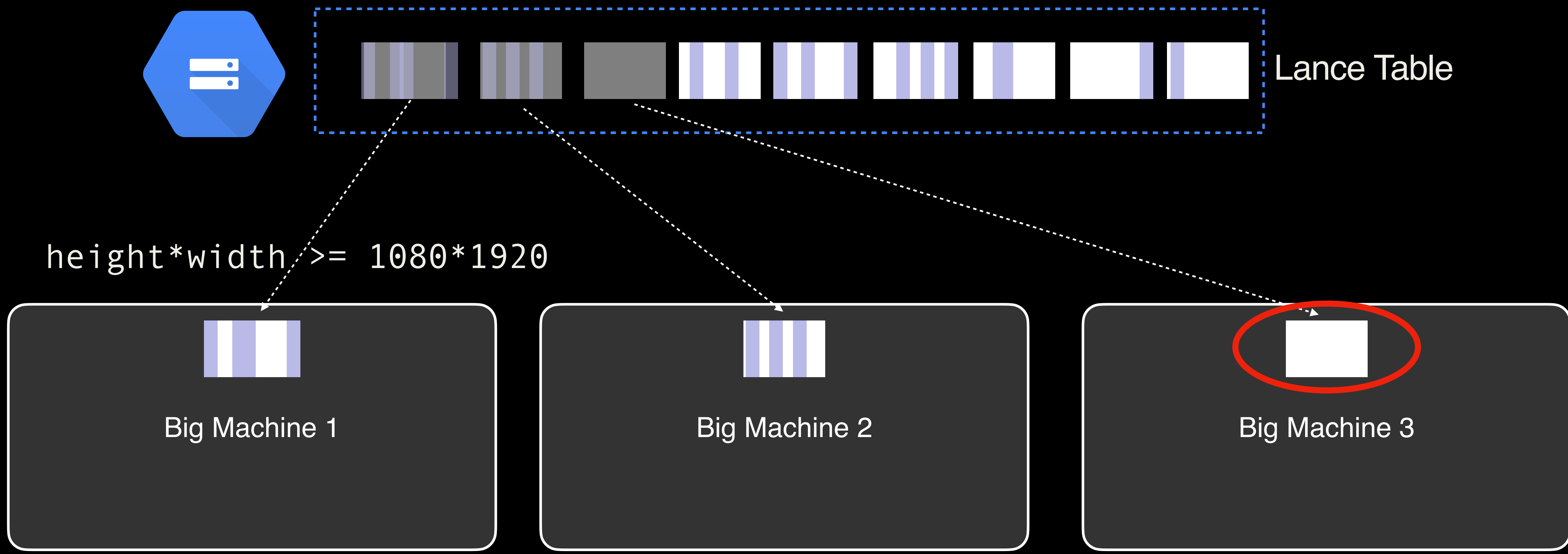
Big Machine 2

Big Machine 3

Lance Fragment Partitioning



Lance Fragment Partitioning

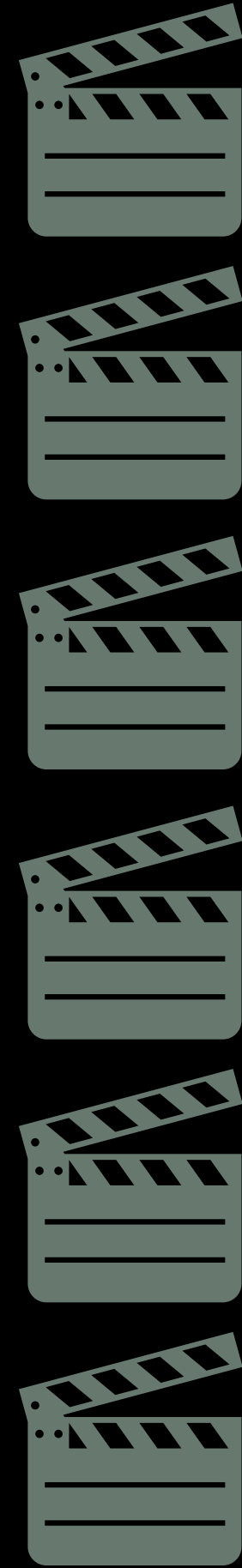




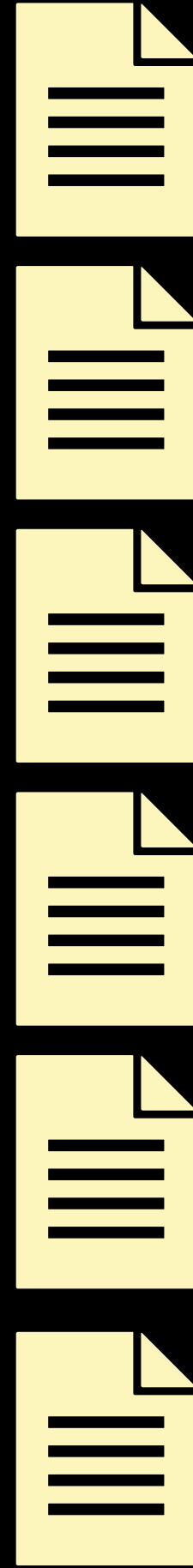
Choose your pointers wisely

POINTERS

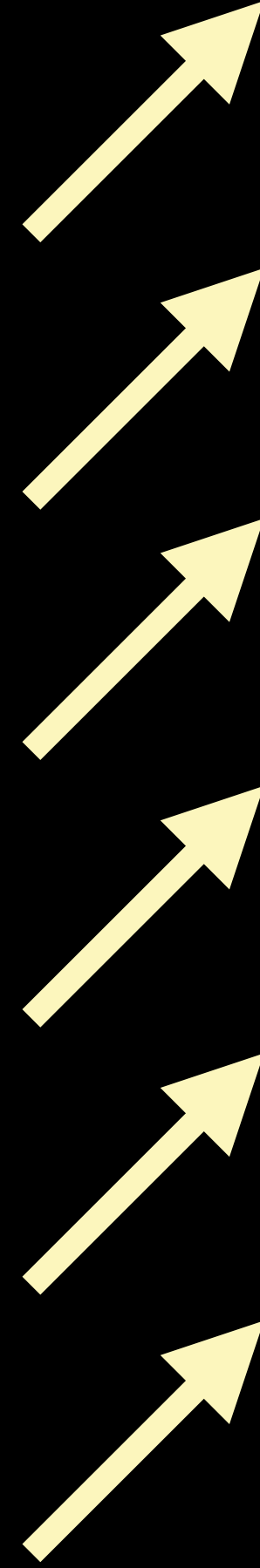
Videos



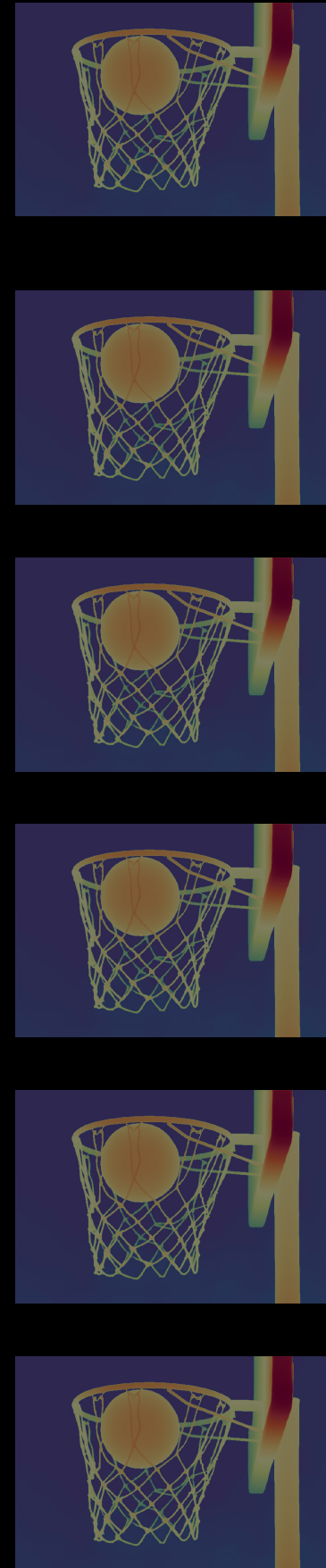
Prompts



Embeddings



Depth Maps



resolution, fps, duration, ...



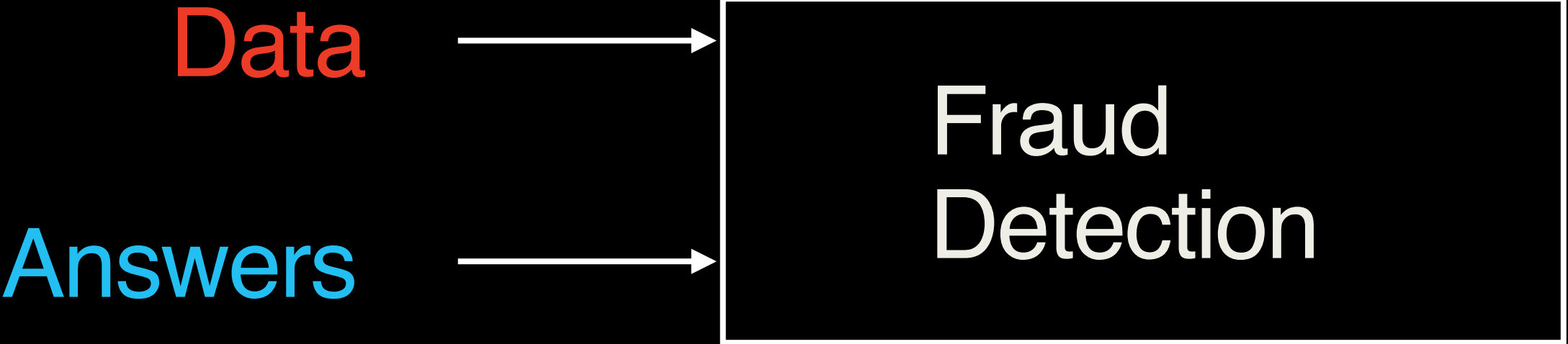
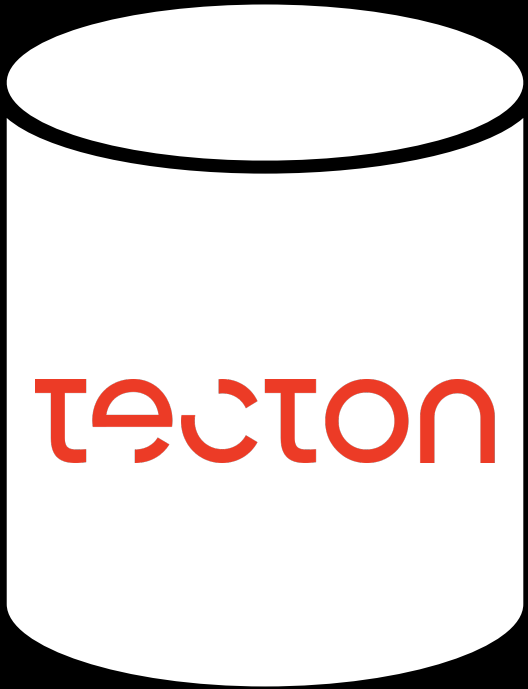
Generative
Video



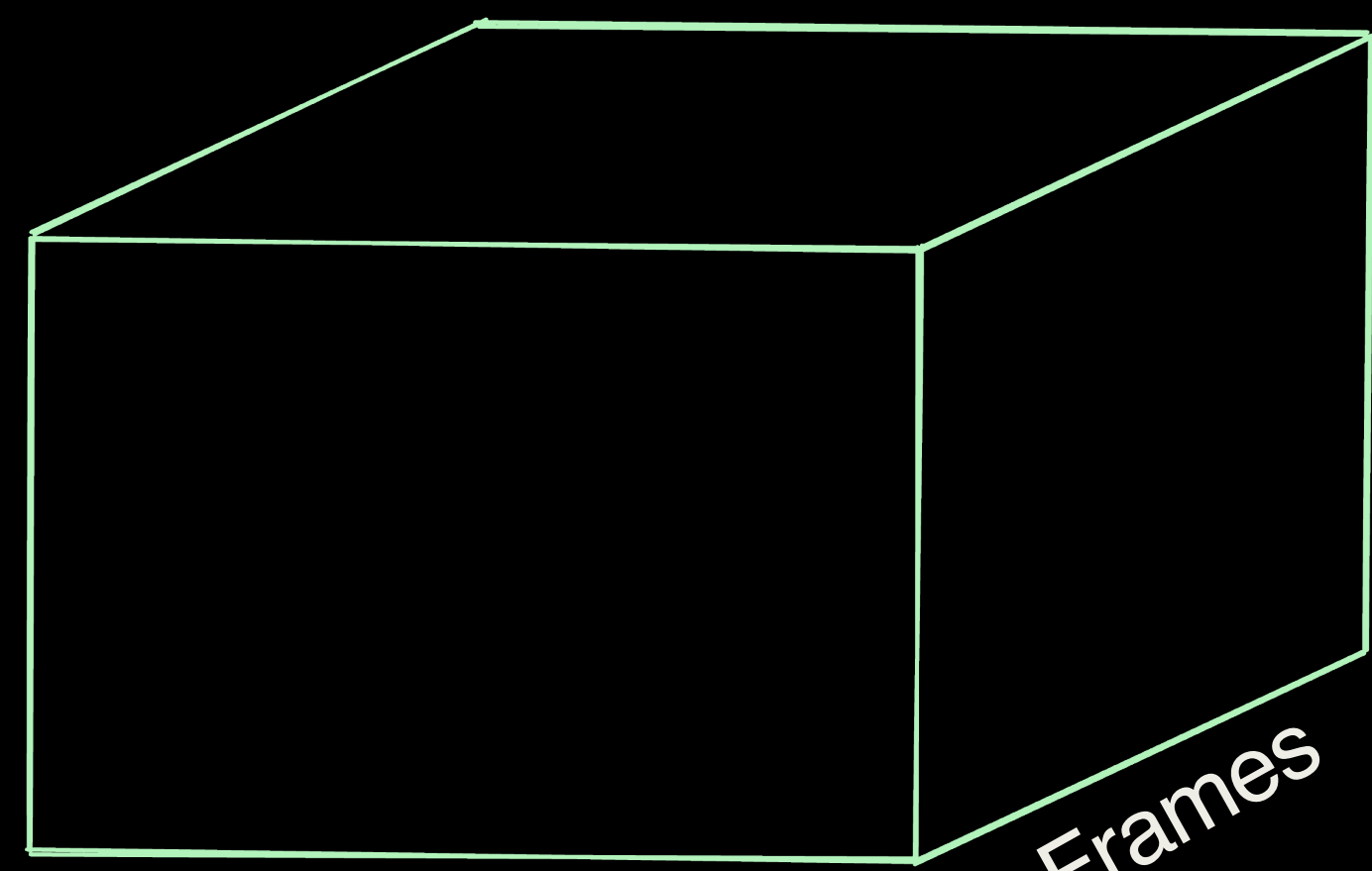
Maximize Read Flexibility



Is Fraud?	IP Address, qps, dollar amt, etc...			



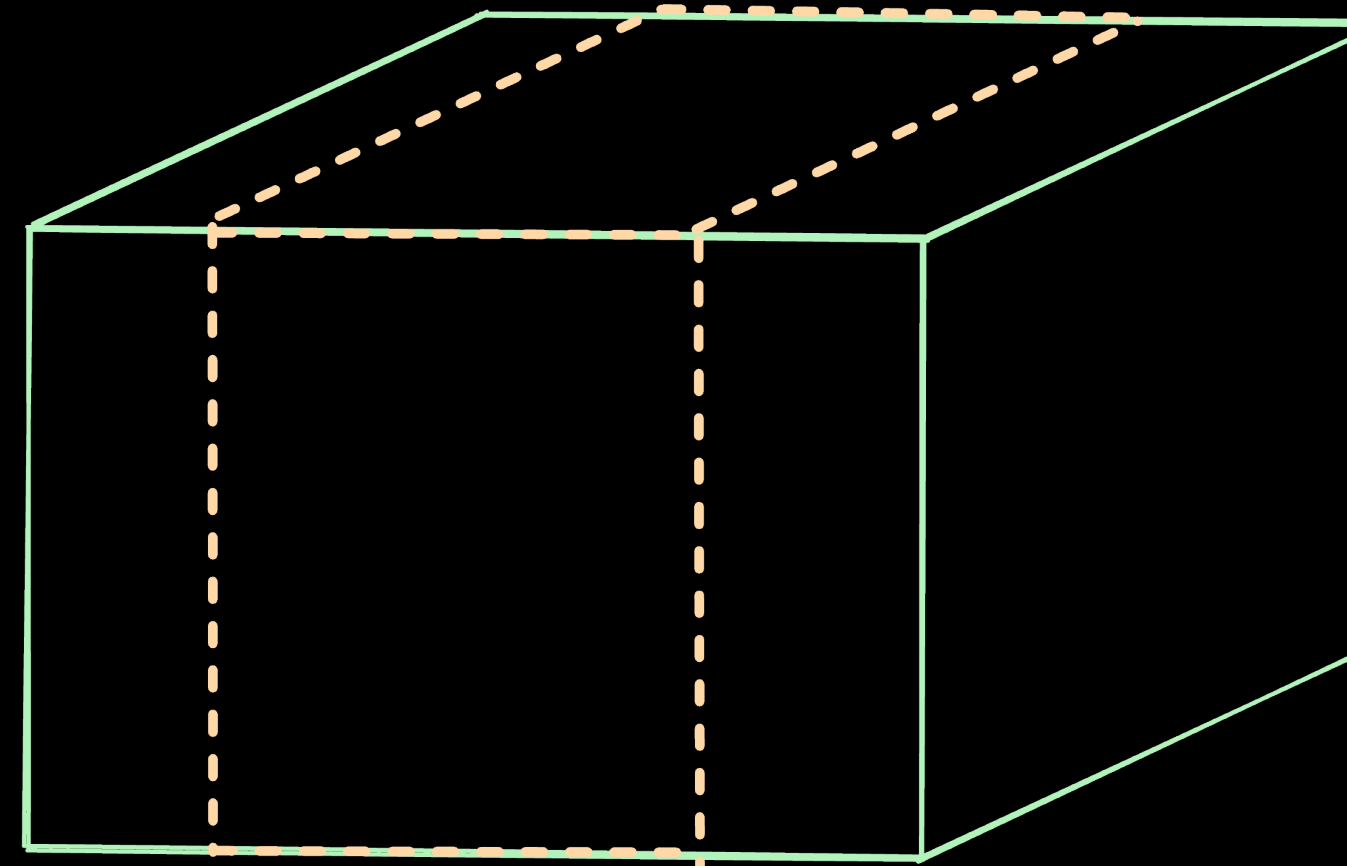
Height



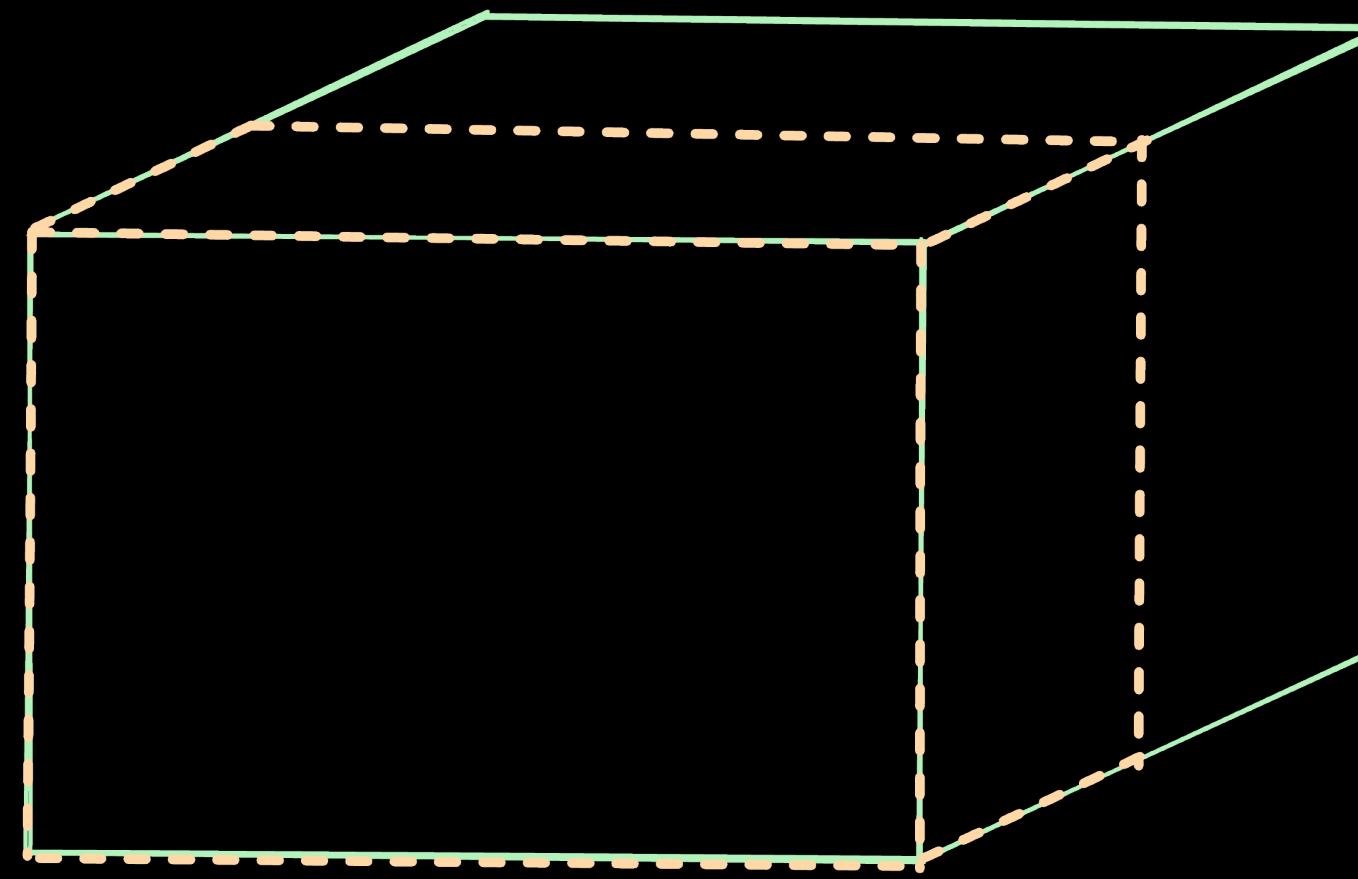
Width

Frames

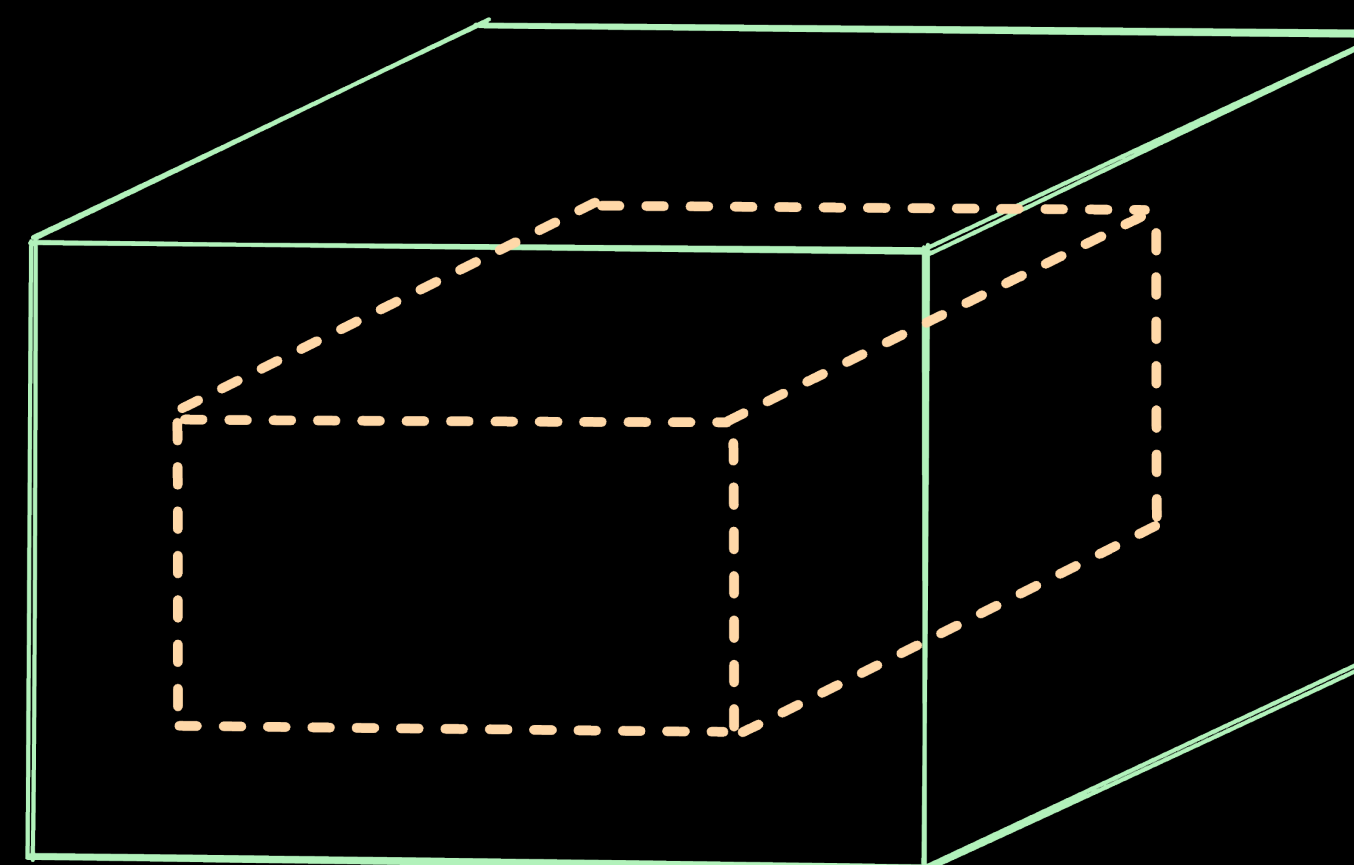
Crop



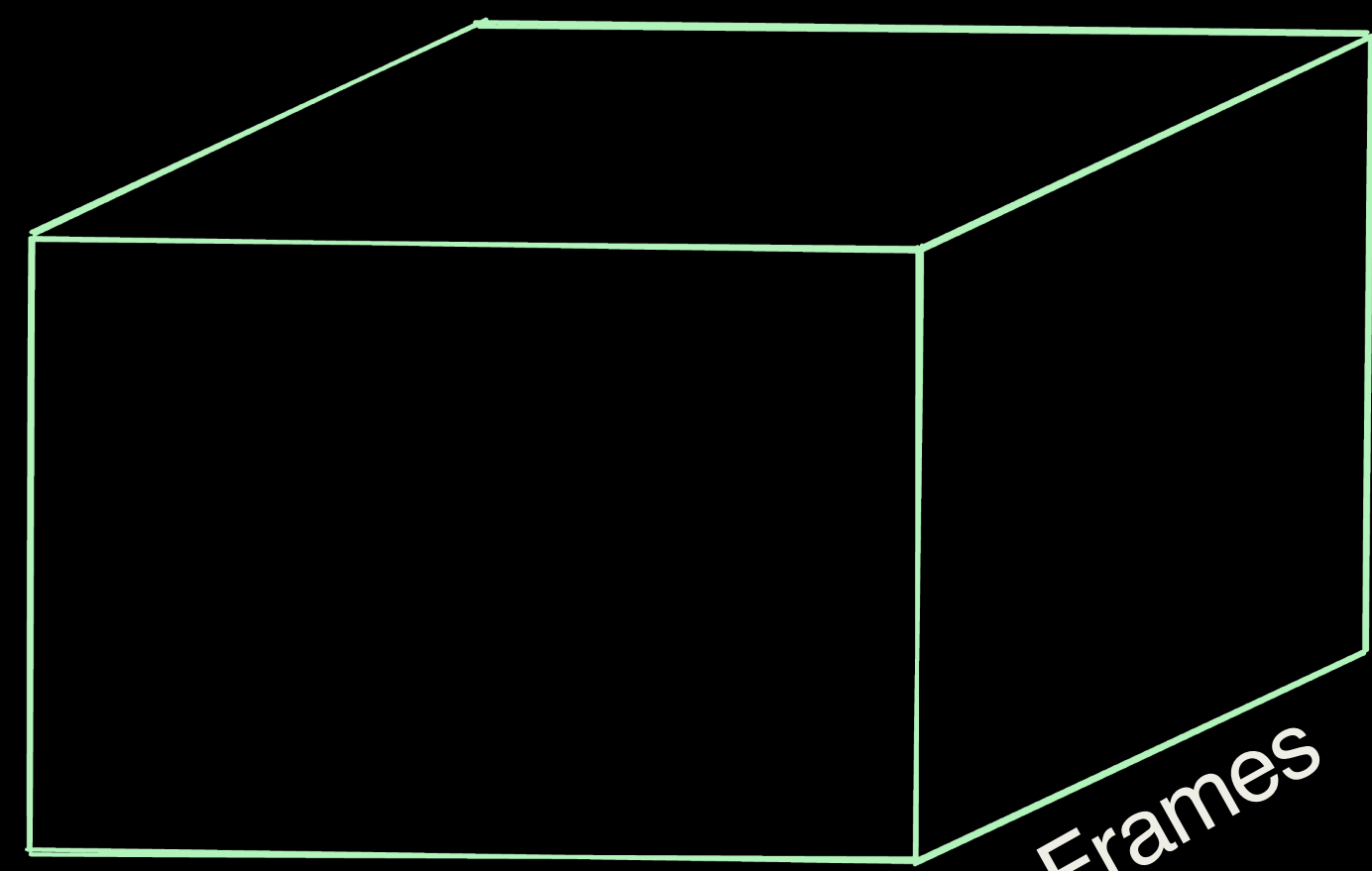
Clip



Resize



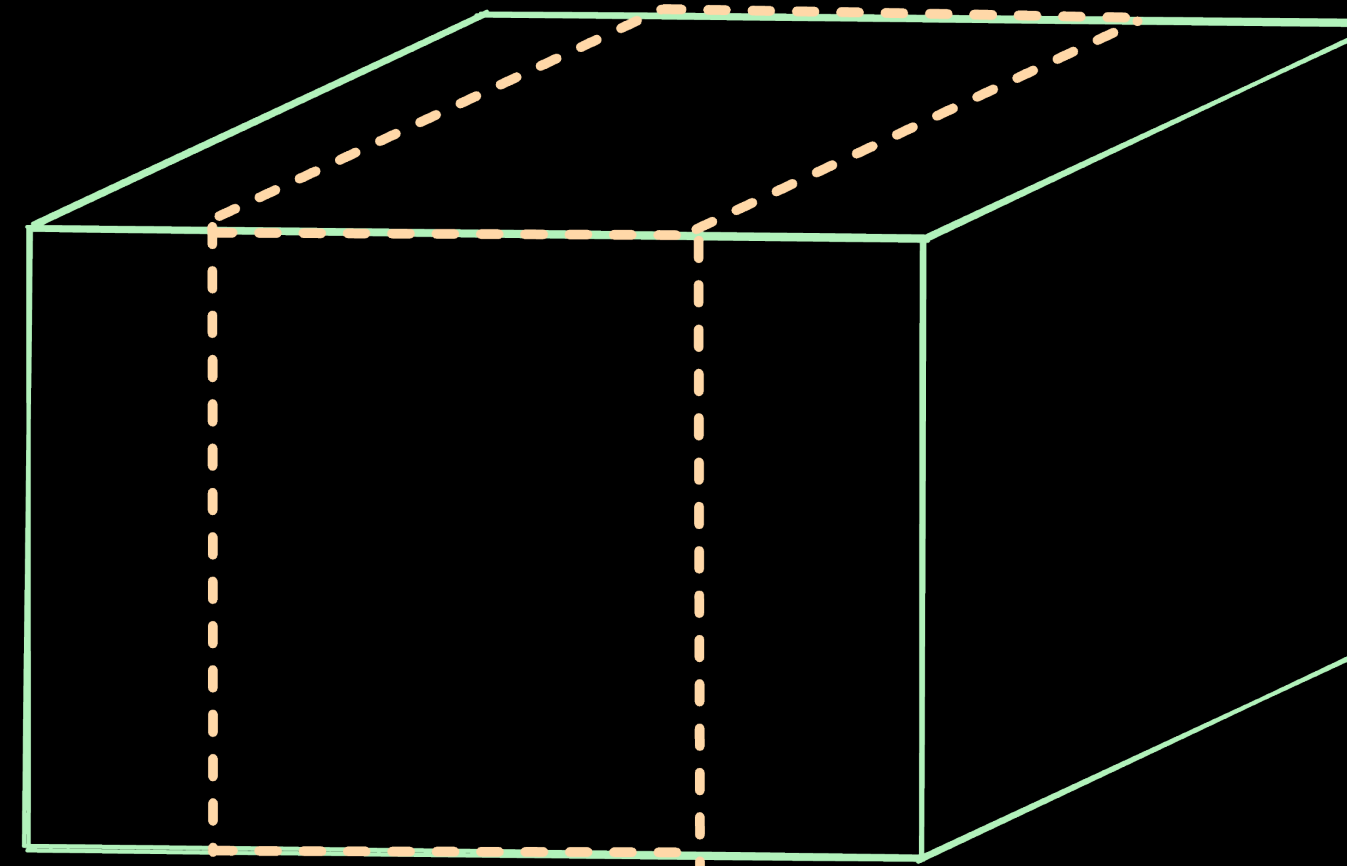
Height



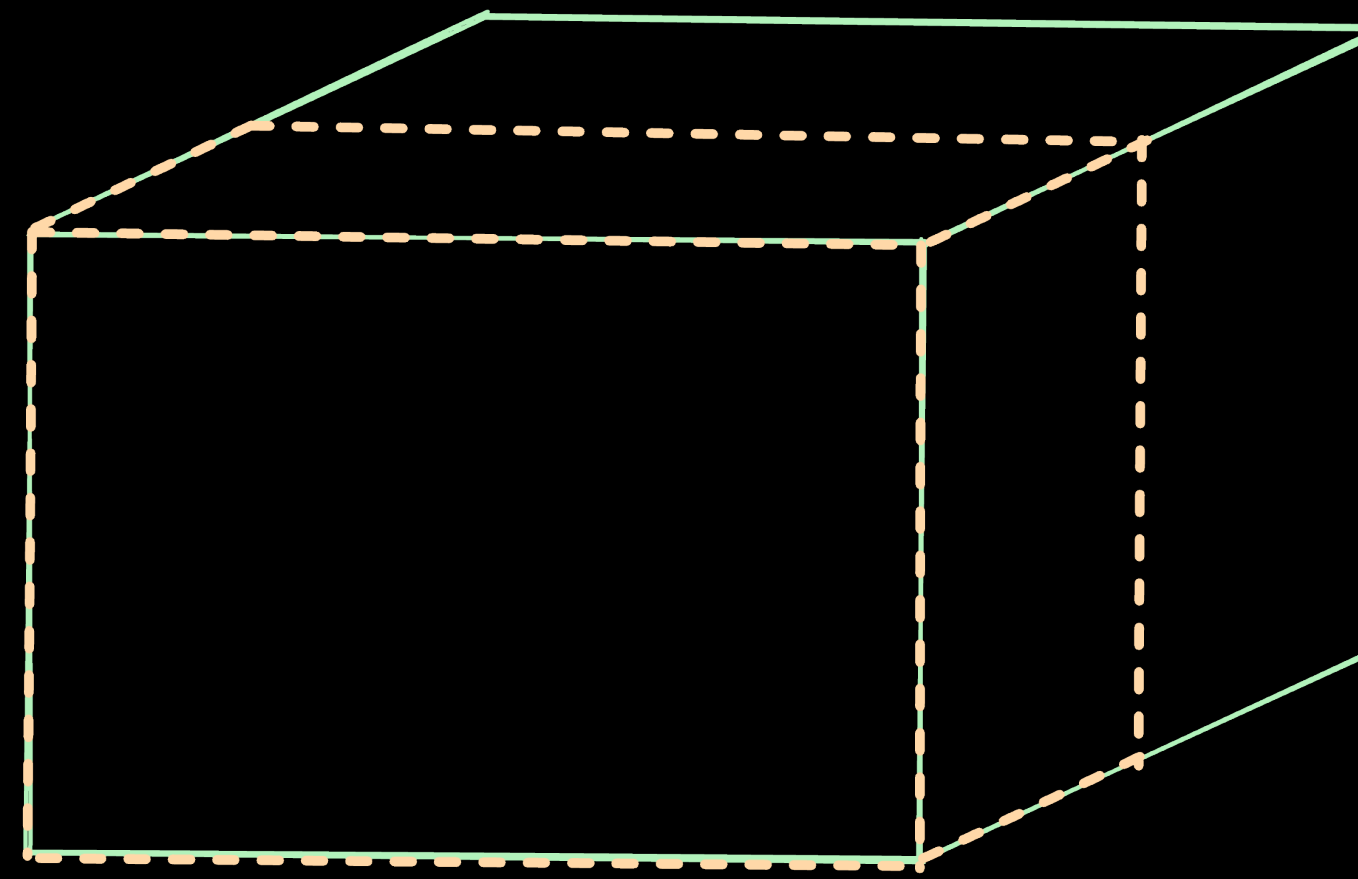
Width

Frames

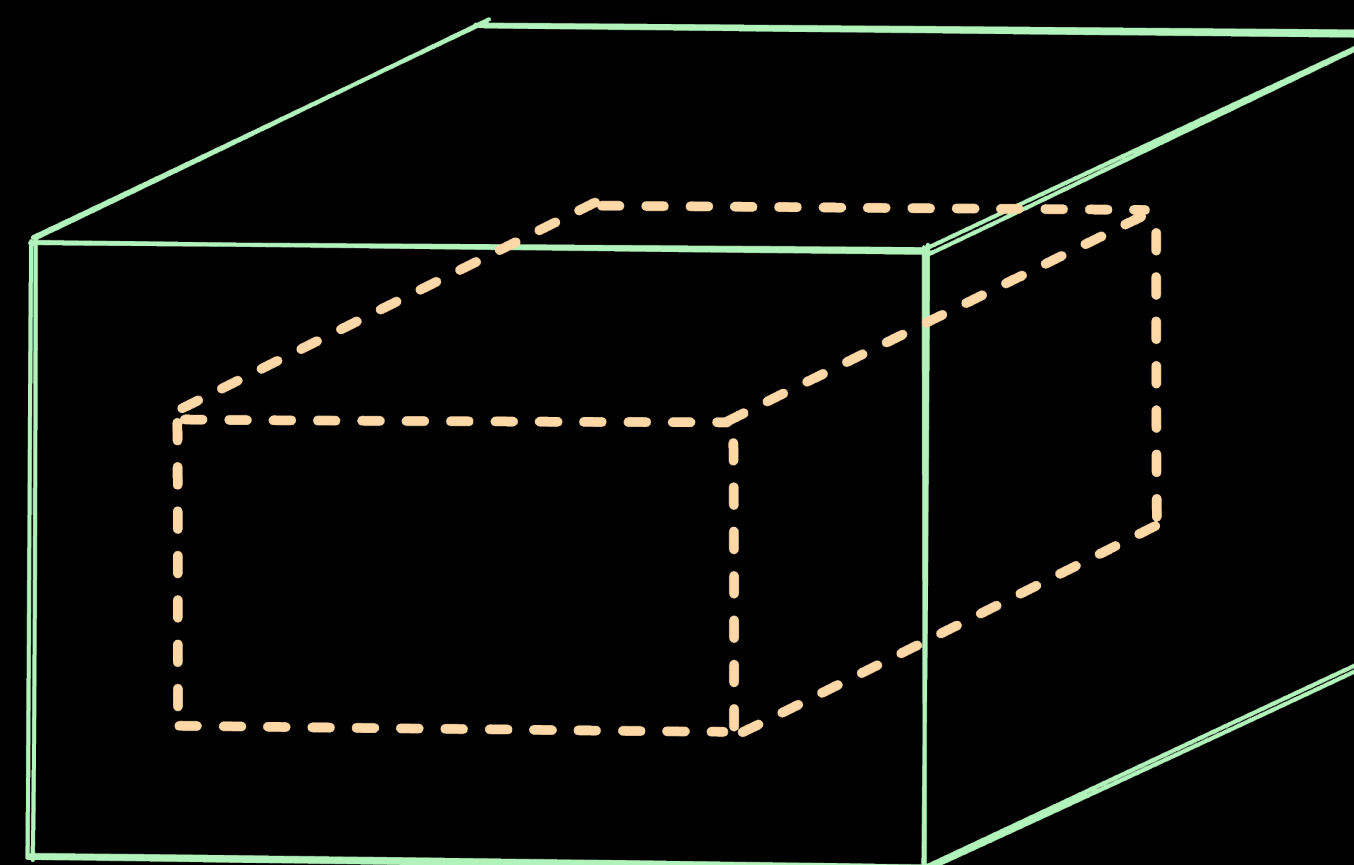
Crop



Clip



Resize



On the Fly

Precompute

rows >> batches

```
for row in rows:
    for processor in processors:
        row = processor(row)

    batcher.consume(row)
    if batcher.full:
        yield batcher.emit_batch()
```

runway

(we're hiring)



ethanrosenthal.com
[@ethanrosenthal.com](https://twitter.com/ethanrosenthal)
[@EthanRosenthal](https://github.com/EthanRosenthal)