# How Beam Serves Models with vLLM
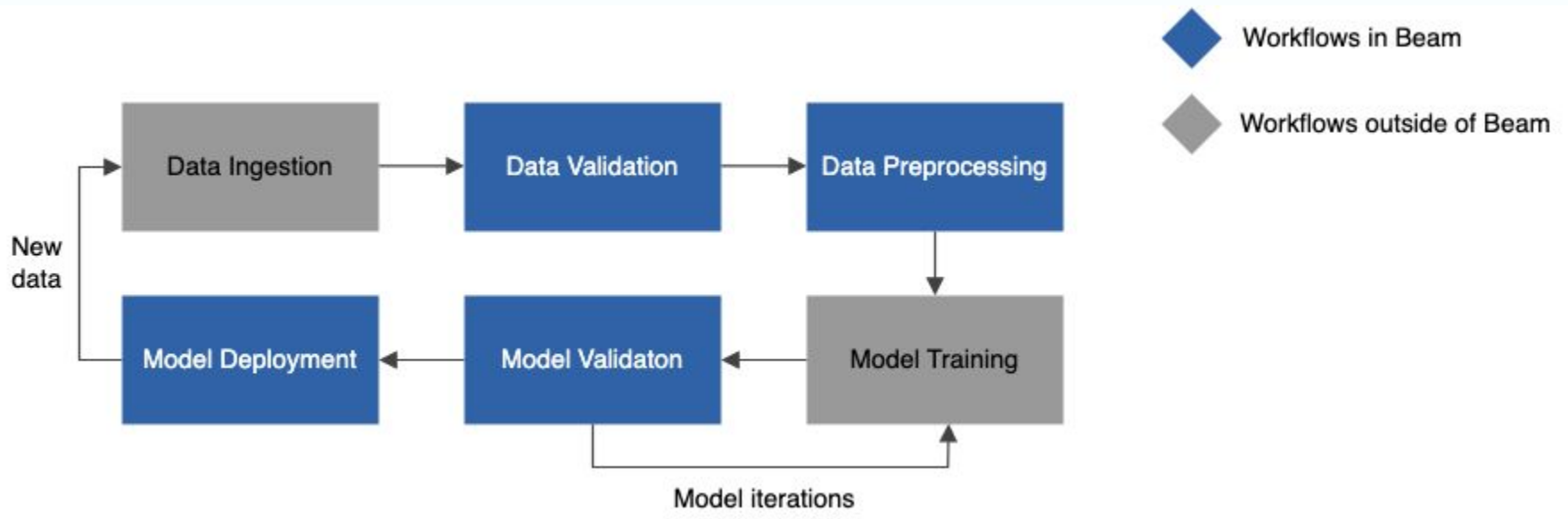
Danny McCormick

# Challenges of Distributed Inference

- Efficiently loading models
- Batching
- Model Updates
- Using multiple models

- Beam takes care of all of this with the RunInference transform
- Loads model, batches inputs, handles updates, and plugs into DAG

```
RunInference(model_handler=<config>)
```

```
>>> data = numpy.array([10, 40, 60, 90],
...                        dtype=numpy.float32).reshape(-1, 1)

>>> model_handler = PytorchModelHandlerTensor(
...     model_class=LinearRegression,
...     model_params={'input_dim': 1, 'output_dim': 1},
...     state_dict_path='gs://path/to/model.pt')

>>> with beam.Pipeline() as p:
...     predictions = (
...         p
...         | beam.Create(data)
...         | beam.Map(torch.Tensor) # Map np array to Tensor
...         | RunInference(model_handler=model_handler)
...         | beam.Map(print))
```

colab.sandbox.google.com/github/apache/beam/blob/master/examples/notebooks/beam-ml/run_inference_huggingface.ipynb

# What is vLLM?

- Open source library for serving large language models
- Takes advantage of the specific architectures of LLMs to perform optimizations a generalized framework can't
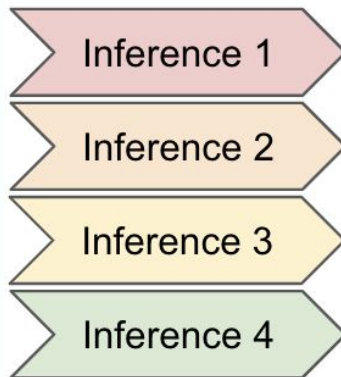  - Examples: continuous batching, PagedAttention

- ML frameworks operate more efficiently when running multiple records in parallel
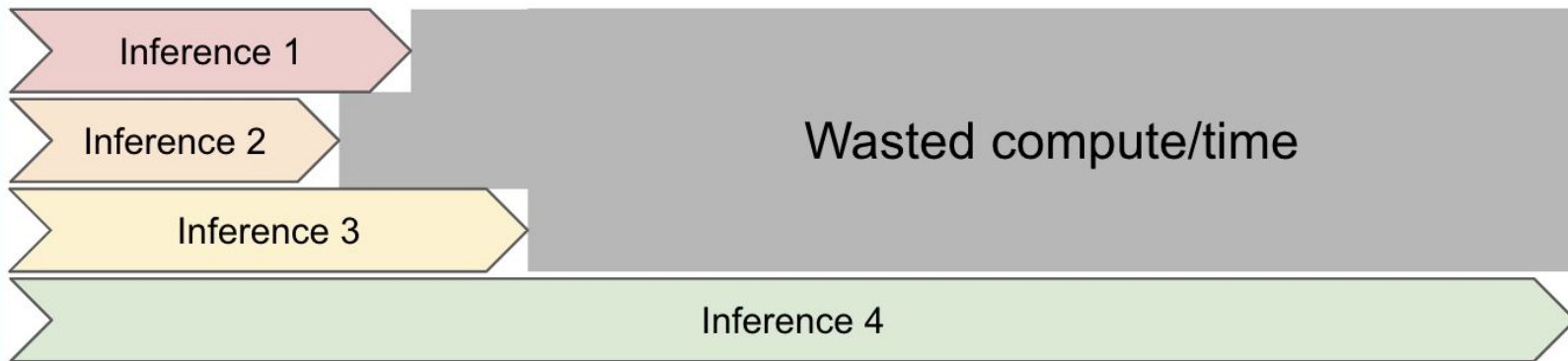


**Without Batching**

Inference 1 → Inference 2 → Inference 3 → Inference 4

**With Batching**

Inference 1
Inference 2
Inference 3
Inference 4

- If some records finish early, they have to wait for the others

- LLMs basically perform an inference (or chain of inferences) per token
- The longer the input, the longer each inference takes (and often the more tokens that need to be generated)
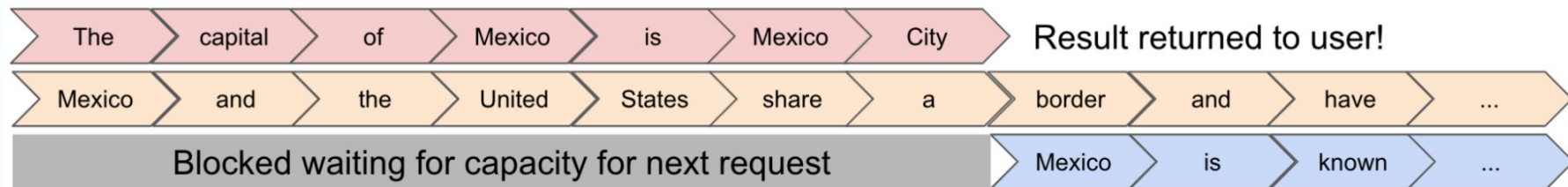
The capital of Mexico is Mexico City | Wasted compute/time

Mexico and the United States share a border and have intertwined histories, but they also have distinct cultural differences. Here are some key similarities and differences...
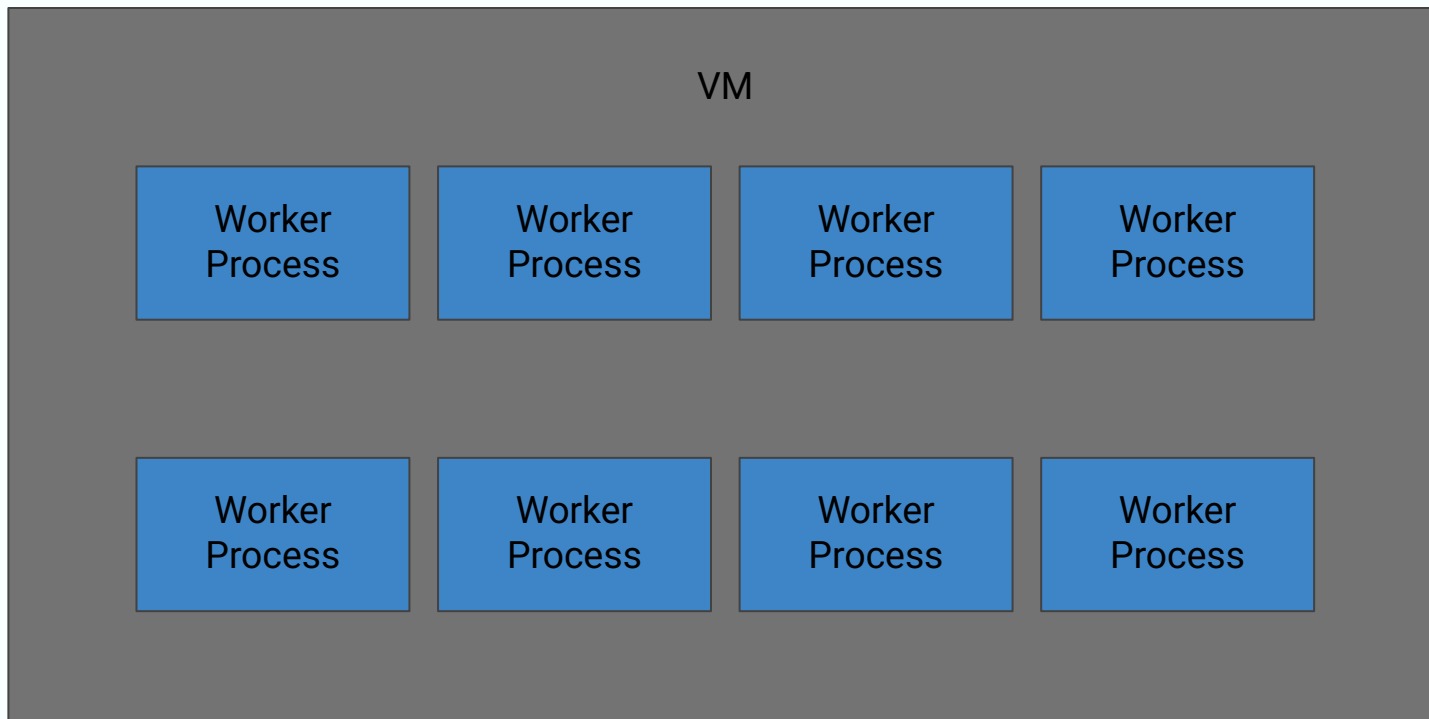
- Batch at the token level instead of the record level

Back to Beam - how Beam does do Model Management

How do we map ideal model configuration to this?
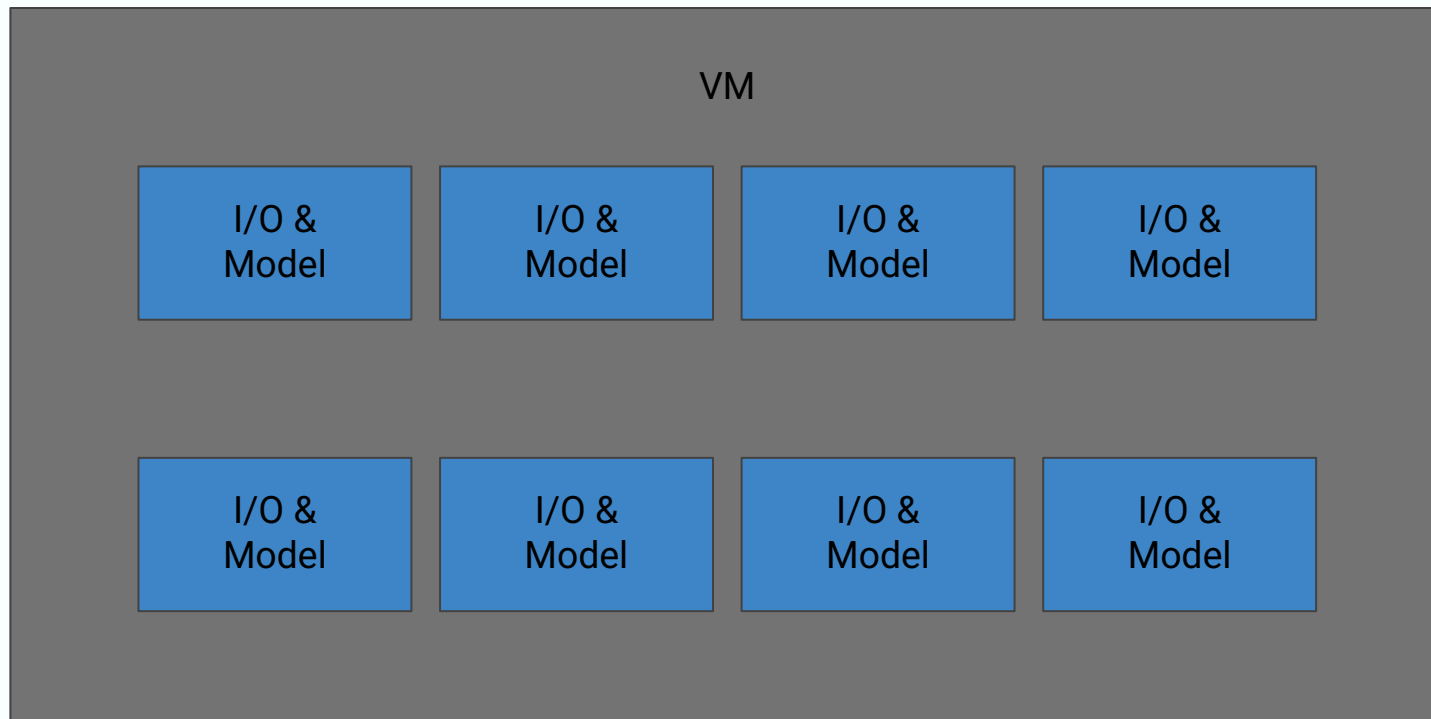
VM

Worker Process (×8)

BEAM SUMMIT NYC 2025
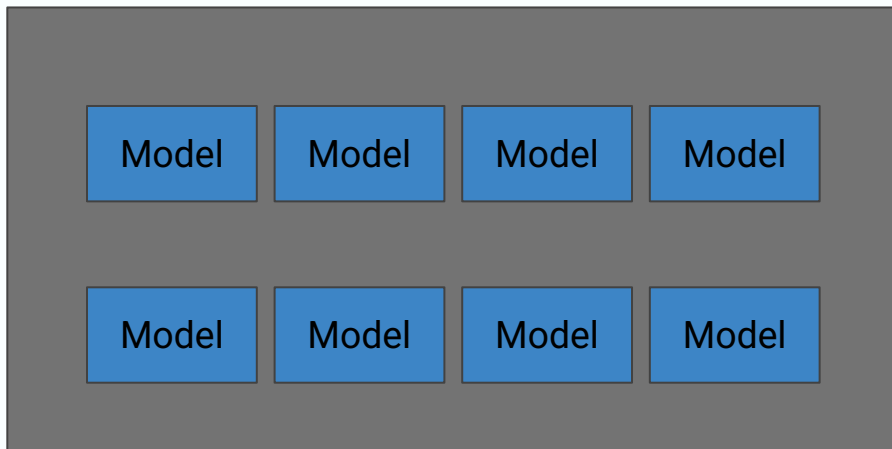
```
>>> model_handler = PytorchModelHandlerTensor(
...     model_class=LinearRegression,
...     model_params={'input_dim': 1, 'output_dim': 1},
...     state_dict_path='gs://path/to/model.pt')

>>> pcoll | RunInference(model_handler=model_handler)
```

```
>>> model_handler = PytorchModelHandlerTensor(
...     model_class=LinearRegression,
...     large_model=True,
...     model_params={'input_dim': 1, 'output_dim': 1},
...     state_dict_path='gs://path/to/model.pt')

>>> pcoll | RunInference(model_handler=model_handler)
```
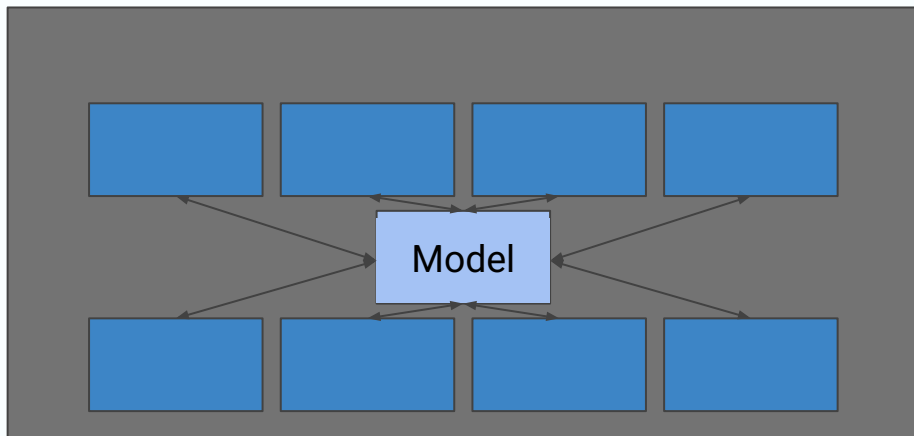
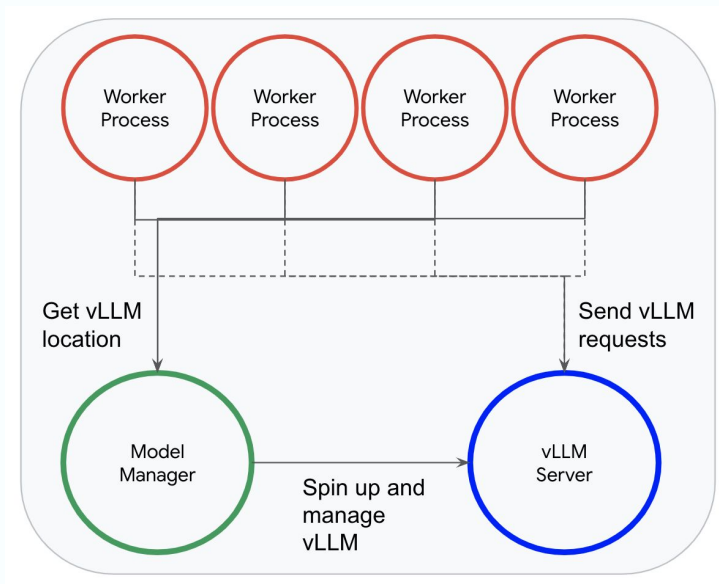- Model Manager (Inference process) now spins up a reference to vLLM server
- Workers talk to vLLM server directly
- Model manager manages vLLM lifecycle
  - Start up
  - Dealing with stuckness/failures
  - Teardown

```
>>> prompts = ["One cause of the console being blank is", "If
you're experiencing network issues","If the button isn't working",
"If you can't submit your job"]

>>> mh = VLLMCompletionsModelHandler('my_favorite_llm')

>>> with beam.Pipeline() as p:
...    predictions = (
...        p
...        | beam.Create(prompts)
...        | RunInference(model_handler=mh)
...        | beam.Map(print))
```
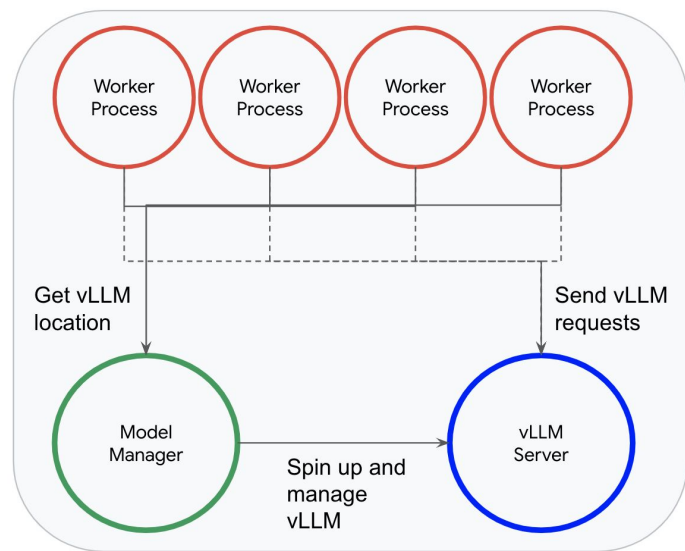
colab.sandbox.google.com/github/apache/beam/blob/master/examples/notebooks/beam-ml/run_inference_vllm.ipynb

- Varies greatly model to model
- With one example pipeline using Google's Gemma 2b model, saw a 23x reduction in number of CPU/GPU core hours when switching from pure Pytorch to vLLM

Danny Mccormick

QUESTIONS?

github.com/damccorm
linkedin.com/in/danny-mccormick-a044b1103

BEAM
SUMMIT
NYC 2025