# How Beam ML Serves Large Models

Danny McCormick
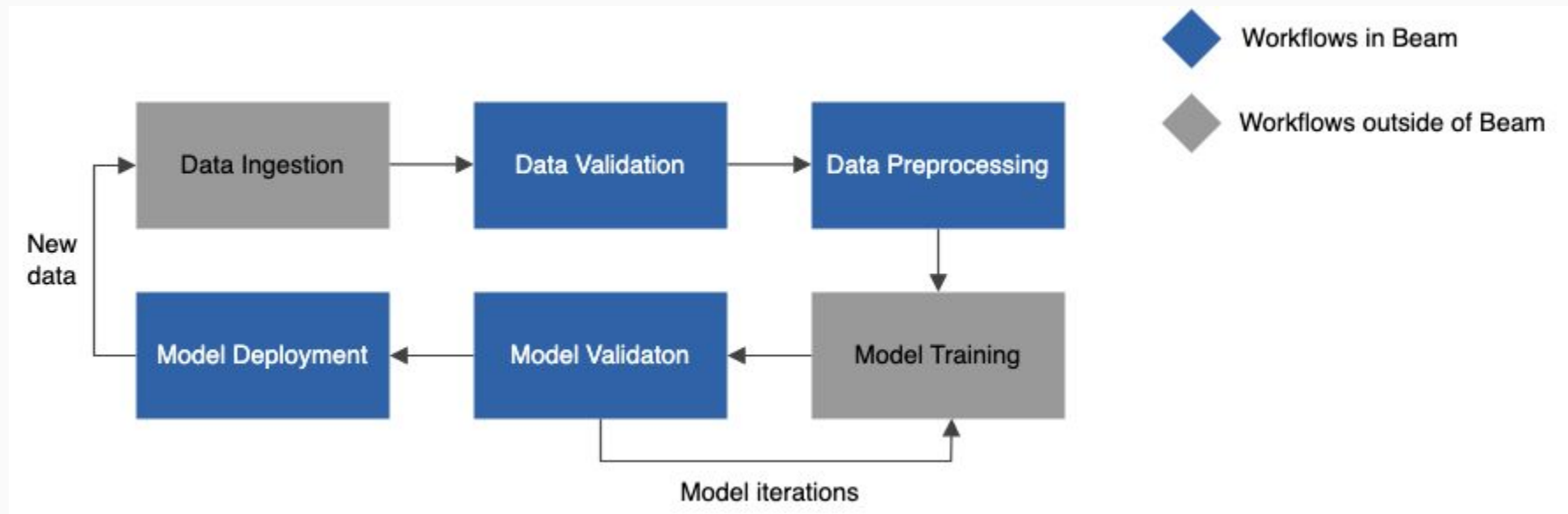
**BEAM SUMMIT**

September 4-5, 2024

Sunnyvale, CA. USA

# The ML lifecycle

# Inference with Beam

- 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))
```
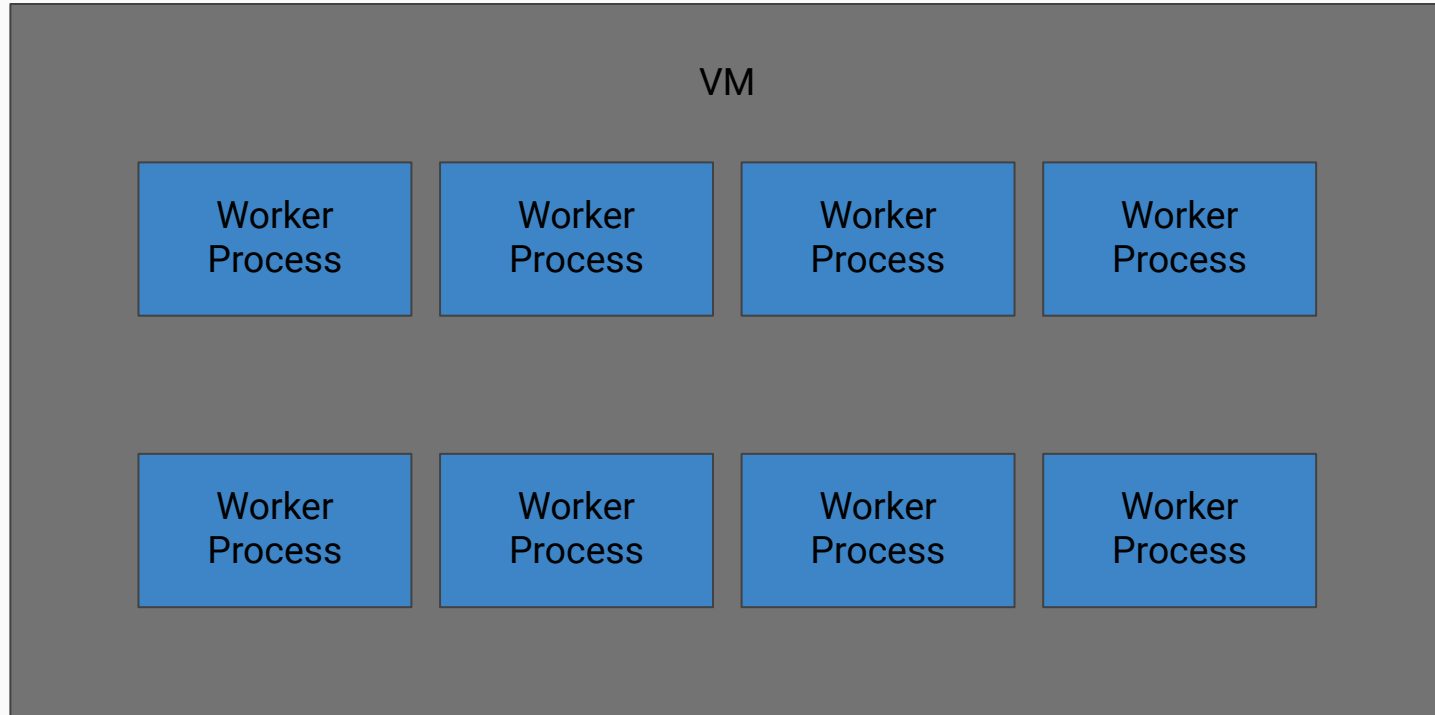
# Basic Inference Demo

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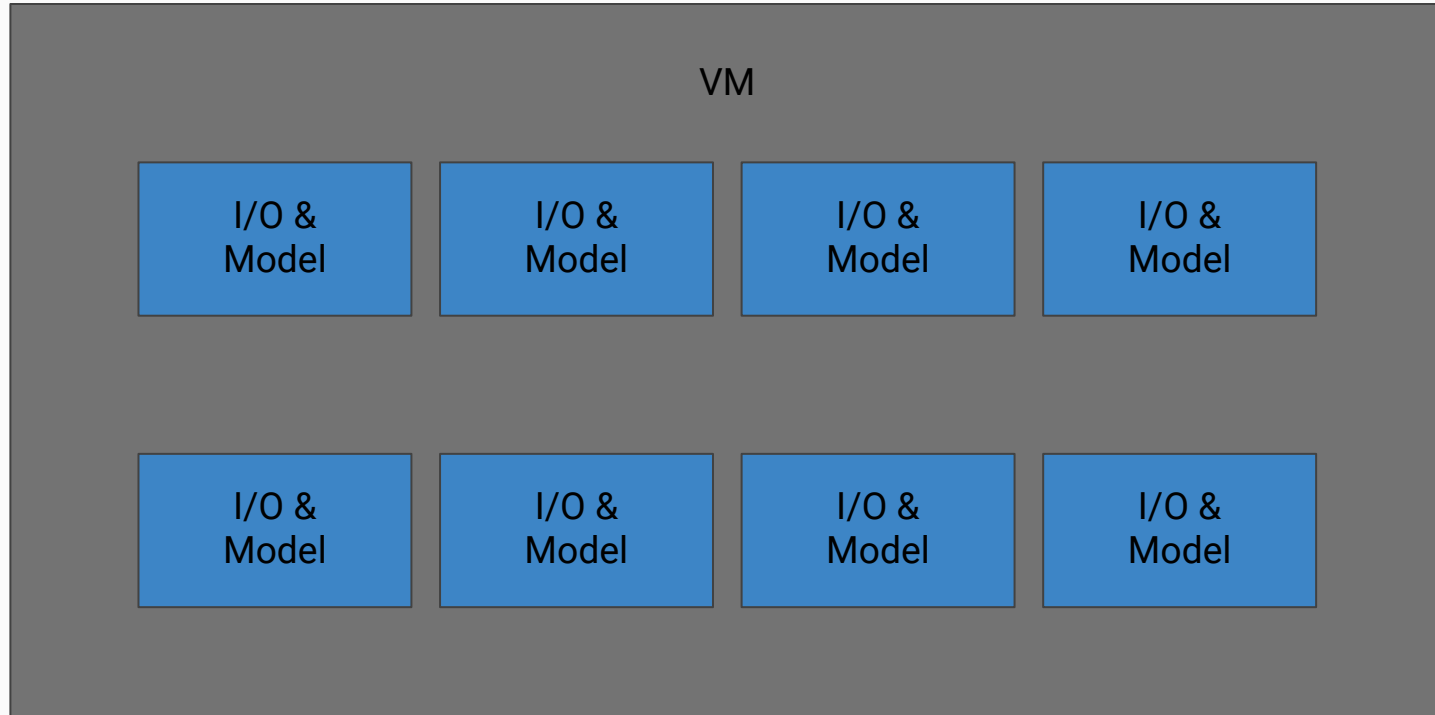
# Large Models

- Text summarization with a small-medium sized LLM
- Can fit many copies of model in memory

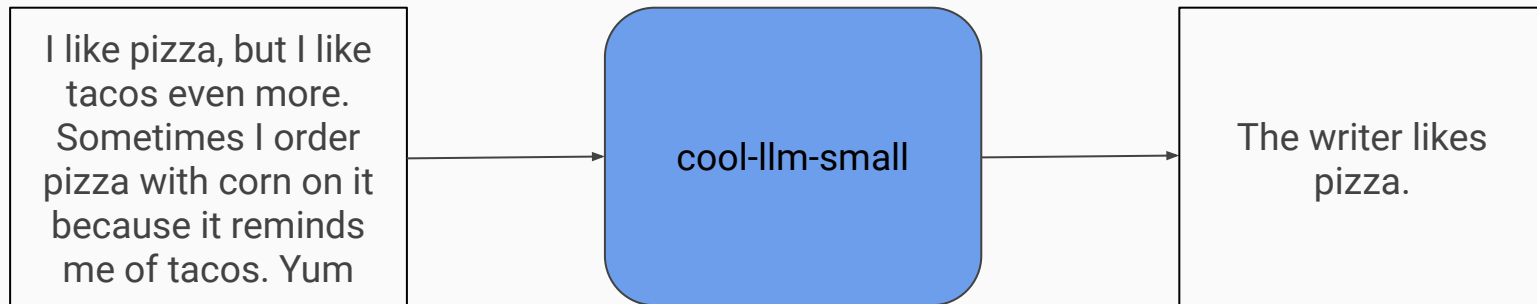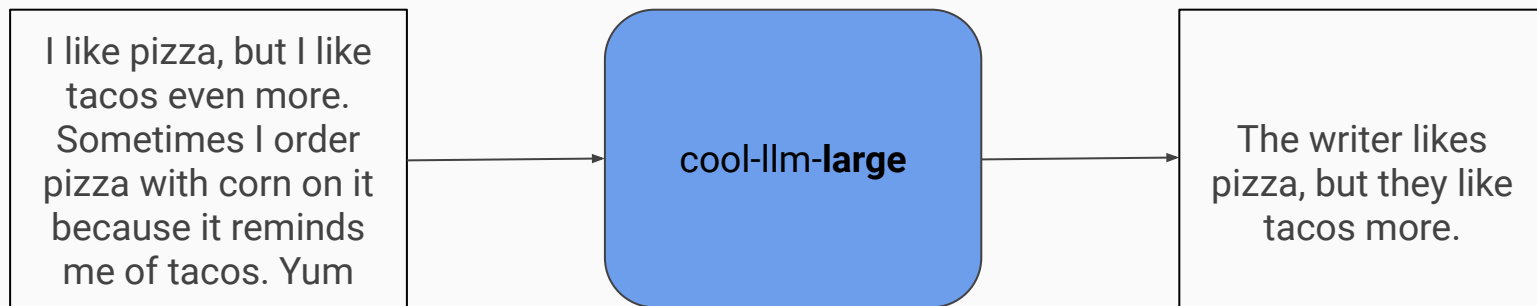| I like pizza so much, it makes me happy. I could eat pizza every day for the rest of my life. | → | cool-llm-small | → | The writer likes pizza. |

# Distributed Runner Architecture*

# Ideal small model configuration

- Aka the easy case
- Uses Beam's shared.py module

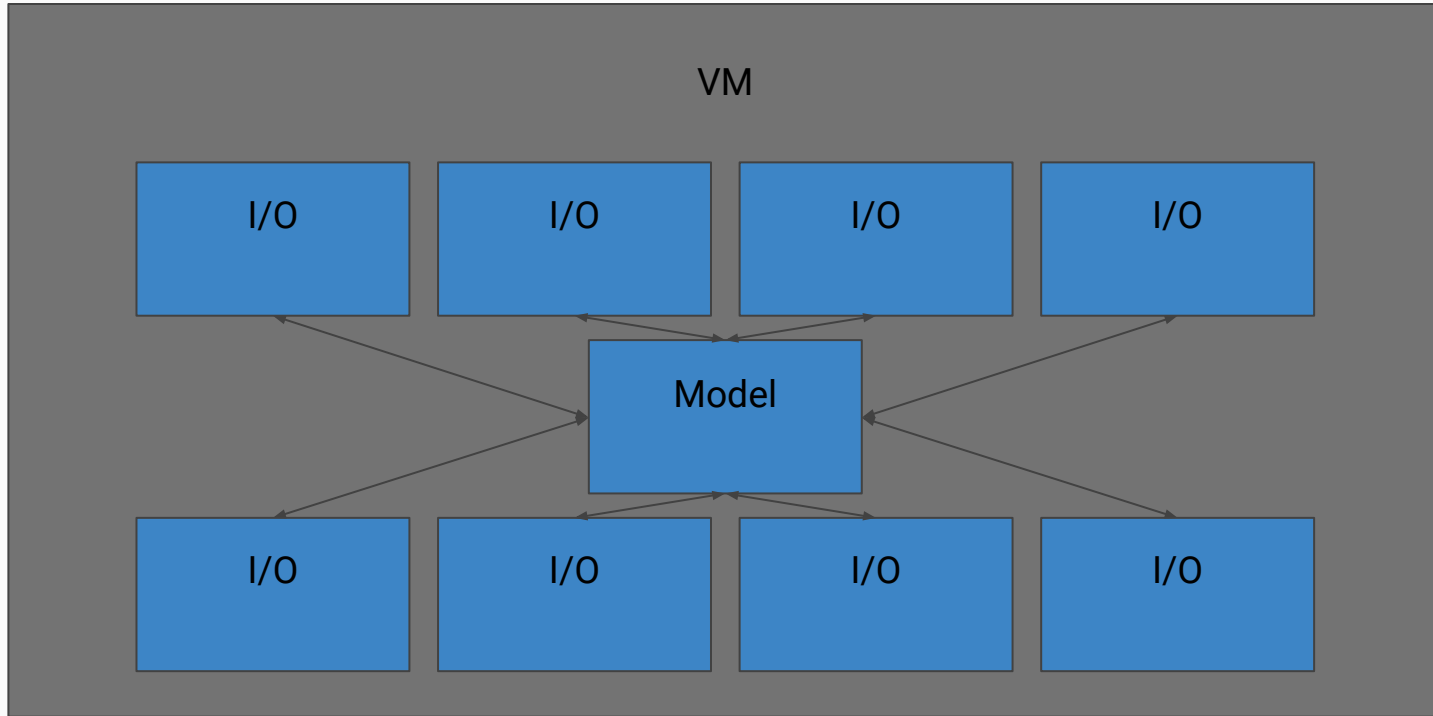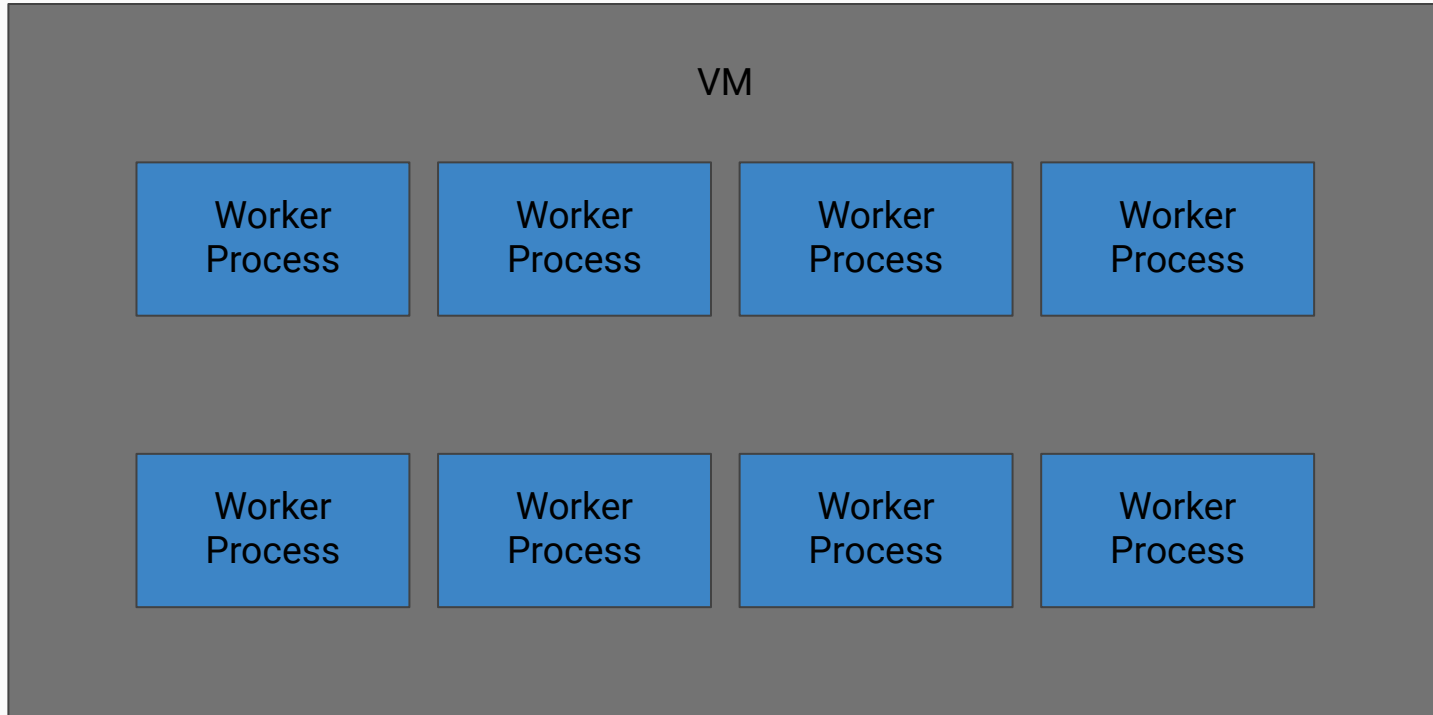- Text summarization with a small-medium sized LLM **isn't good enough**

| I like pizza, but I like tacos even more. Sometimes I order pizza with corn on it because it reminds me of tacos. Yum | → | cool-llm-small | → | The writer likes pizza. |
| --- | --- | --- | --- | --- |

- Many options, one is to switch to a larger model
- Can only fit one (or few) copies in memory

| I like pizza, but I like tacos even more. Sometimes I order pizza with corn on it because it reminds me of tacos. Yum | → | cool-llm-**large** | → | The writer likes pizza, but they like tacos more. |

# Ideal Large Model Configuration

VM

I/O

I/O

I/O

I/O

Model

I/O

I/O

I/O

I/O

# How do we map ideal model configurations to this?

- Reduce memory at cost of interprocess communication, minimized parallelism
- Uses beam's multi_process_shared library

- If you're spinning this yourself, you need to set up a new serving topology, but Beam can make it easy
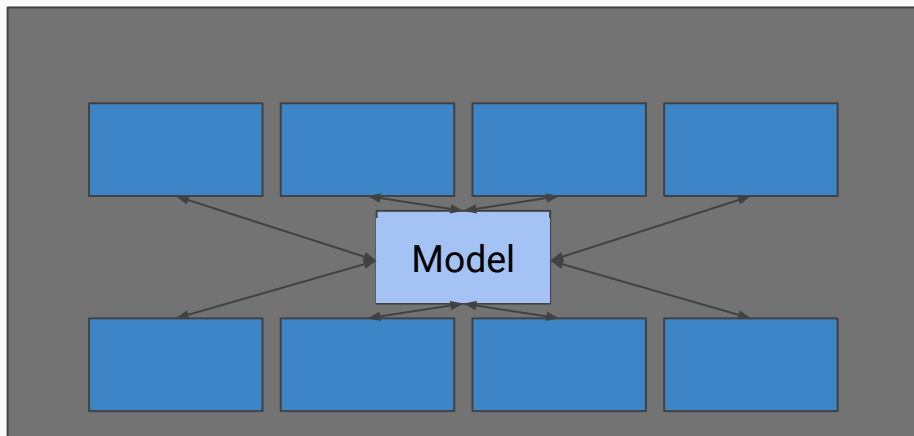
# Built in model handler using default configuration

```
>>> 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)
```

# Built in model handler using large model configuration

```python
>>> 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)
```

# Custom Model Handler configuration (default, share across threads)

```
>>> def run_inference(model, batch, …):
...    model.predict(batch)
```

# Custom Model Handler configuration (large model configuration)

```
>>> def run_inference(model, batch, …):
...    model.predict(batch)

>>> def share_model_across_processes(self) -> bool:
...    return true
```
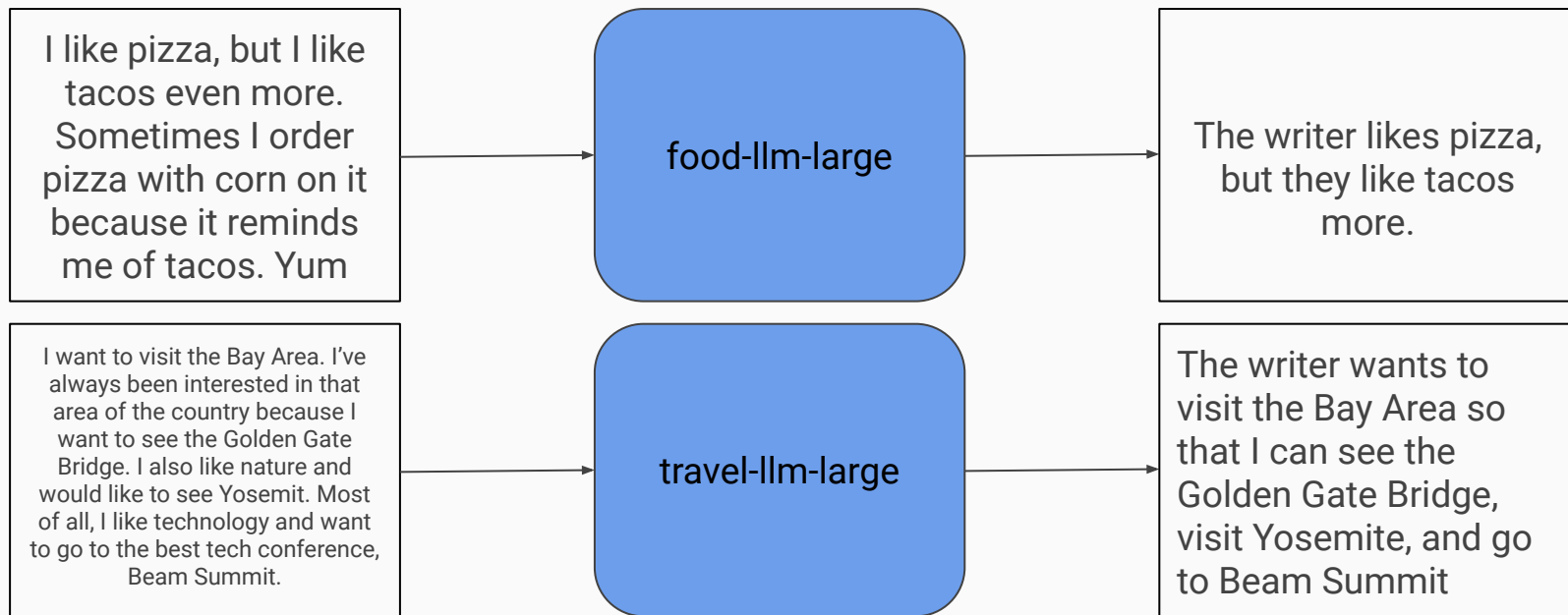
# Custom Model Handler configuration (medium model configuration)

```
>>> def run_inference(model, batch, …):
...    model.predict(batch)

>>> def share_model_across_processes(self) -> bool:
...    return true

>>> def model_copies(self) -> int:
...    return 4
```
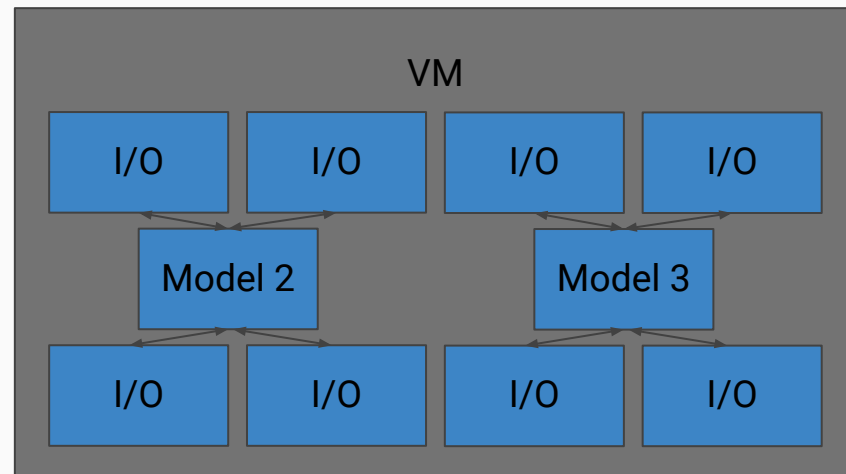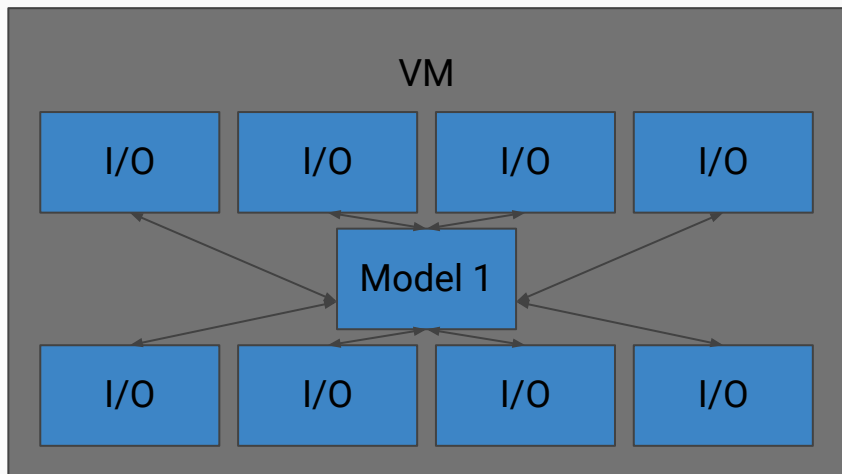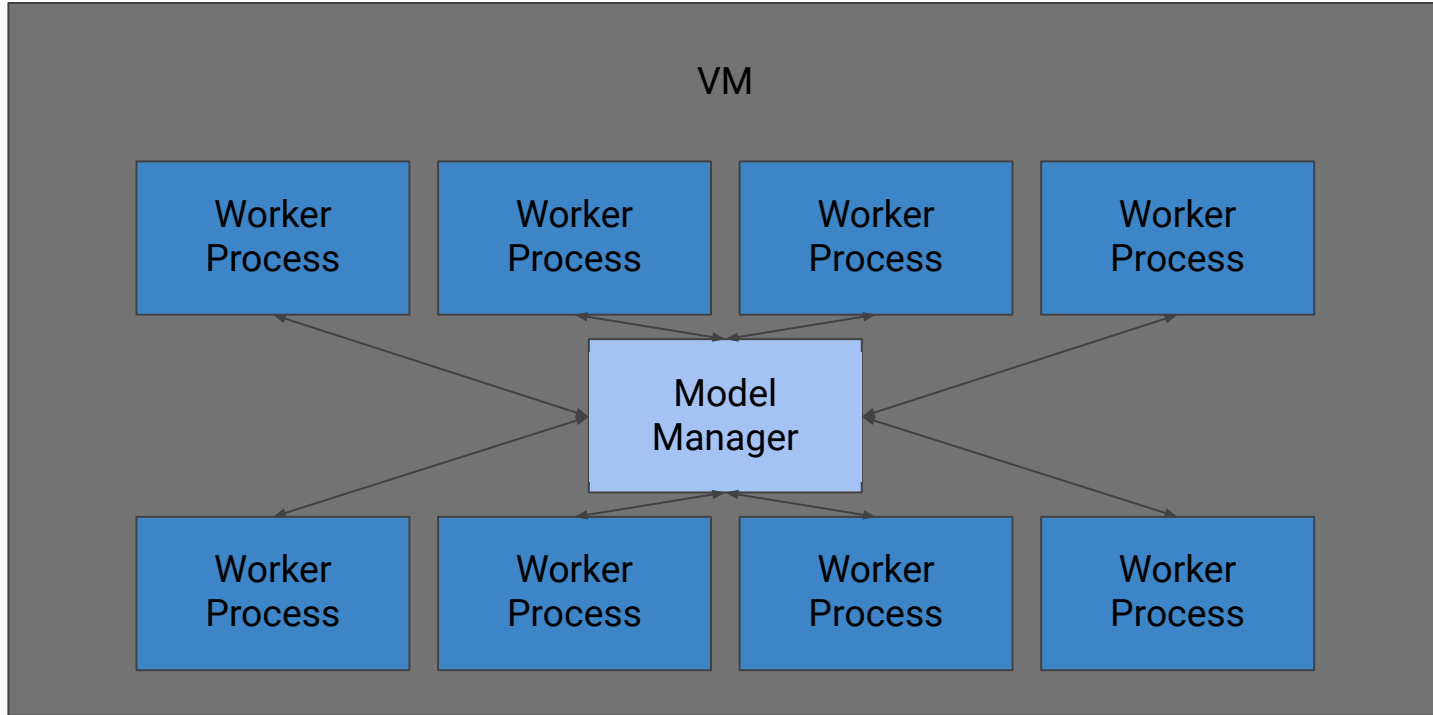
- ## What if I need a model per customer?

| | | |
|---|---|---|
| I like pizza, but I like tacos even more. Sometimes I order pizza with corn on it because it reminds me of tacos. Yum | food-llm-large | The writer likes pizza, but they like tacos more. |
| I want to visit the Bay Area. I've always been interested in that area of the country because I want to see the Golden Gate Bridge. I also like nature and would like to see Yosemit. Most of all, I like technology and want to go to the best tech conference, Beam Summit. | travel-llm-large | The writer wants to visit the Bay Area so that I can see the Golden Gate Bridge, visit Yosemite, and go to Beam Summit |

# Ideal Large Model Configuration

- Model Manager empowered to load/unload models in order to make optimal use of memory

# Again, this is just configuration

```
>>> mh1 = PytorchModelHandlerTensor(
...     model_class=LinearRegression,
...     model_params={'input_dim': 1, 'output_dim': 1},
...     state_dict_path='gs://path/to/model.pt')




>>> pcoll | RunInference(mh1)
```

# Again, this is just configuration

```
>>> mh1 = PytorchModelHandlerTensor(
...     model_class=LinearRegression,
...     model_params={'input_dim': 1, 'output_dim': 1},
...     state_dict_path='gs://path/to/model.pt')
>>> mh2 = <...>

>>> per_key_mhs = [
...     KeyModelMapping(['key1', 'key2', 'key3'], mh1),
...     KeyModelMapping(['foo', 'bar', 'baz'], mh2)]
>>> mh = KeyedModelHandler(per_key_mhs)

>>> pcoll | RunInference(mh)
```
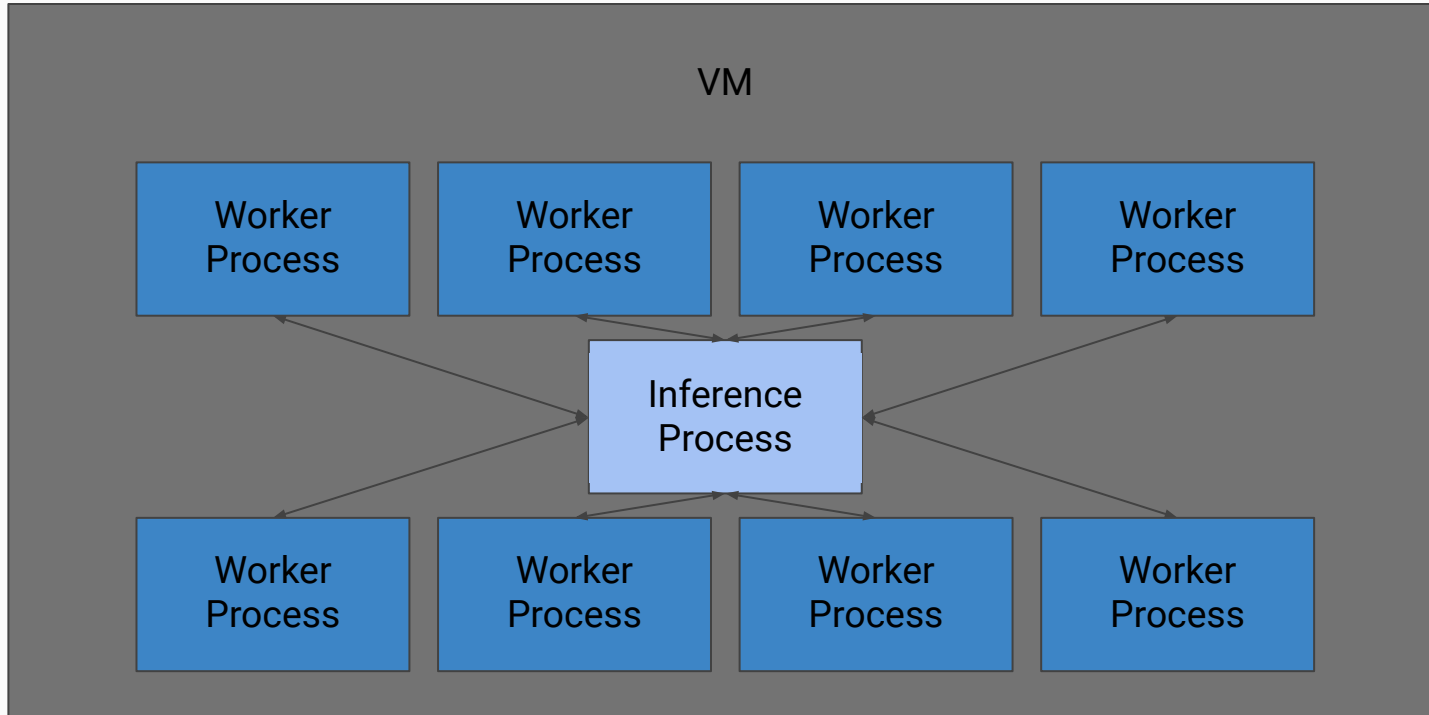
# Per Key Model Demo

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

# Specialty Hardware

- Hardware availability dependent on runner
- Beam has some primitives that help

- Resource hints for heterogeneous pools
- Built in detection + framework specific responses to GPUs at the ModelHandler level
- Large model setting helps

Central Inference Process provides a single point of interaction with GPU

# Try it yourself

https://github.com/apache/beam/tree/master/examples/notebooks/beam-ml

# Thank you!

Questions?

dannymccormick@google.com

Github - damccorm

https://www.linkedin.com/in/danny-mccormick-a044b1103/