# PyTorch/XLA Data Parallel with SPMD
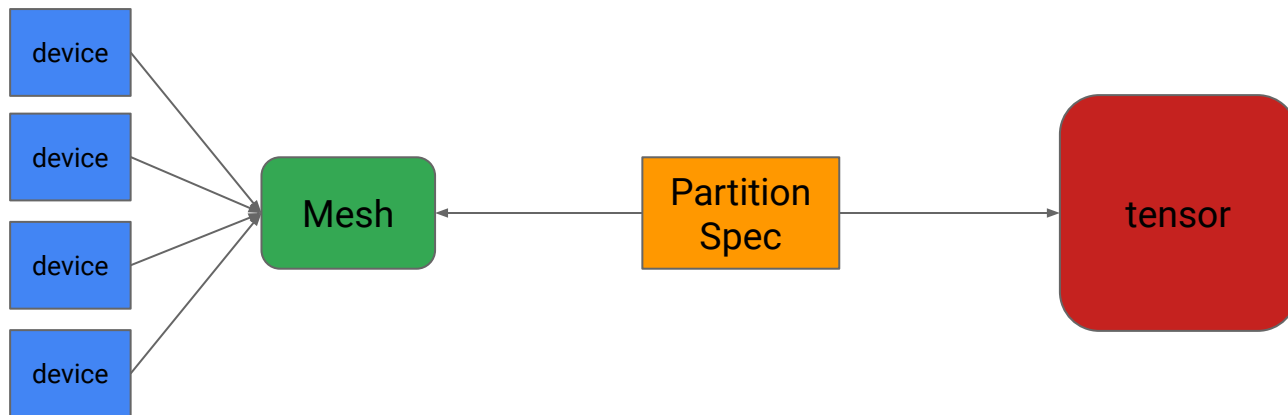
# Agenda

1. Recap of SPMD
2. What's data parallel
3. How data parallel + SPMD works
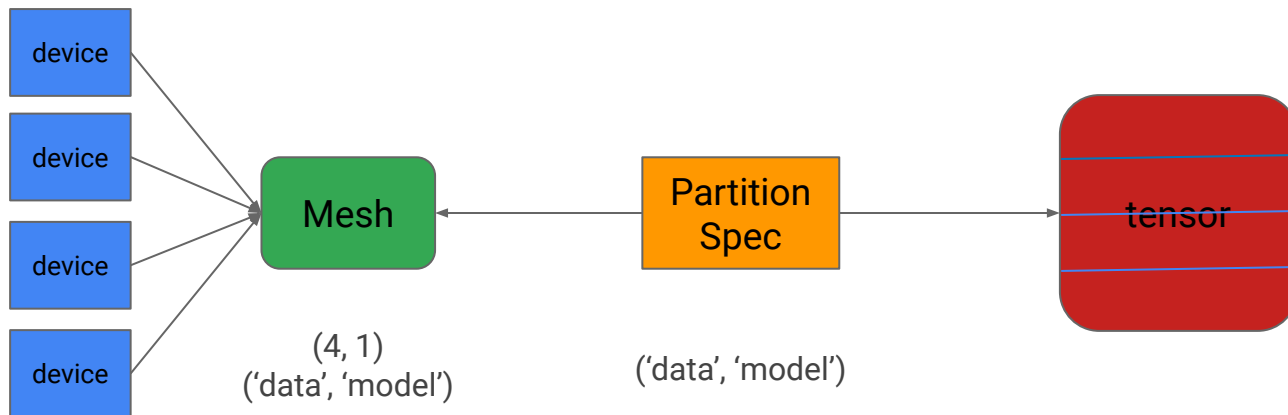
# GSPMD

1. https://arxiv.org/abs/2105.04663
2. User only express sharding intention, let compiler shard the tensor for you.
3. User don't need to shard every tensor, compiler will propagate the sharding for the user.
4. Collective ops(all_gather, reduce_scatter etc) will be added after compilation
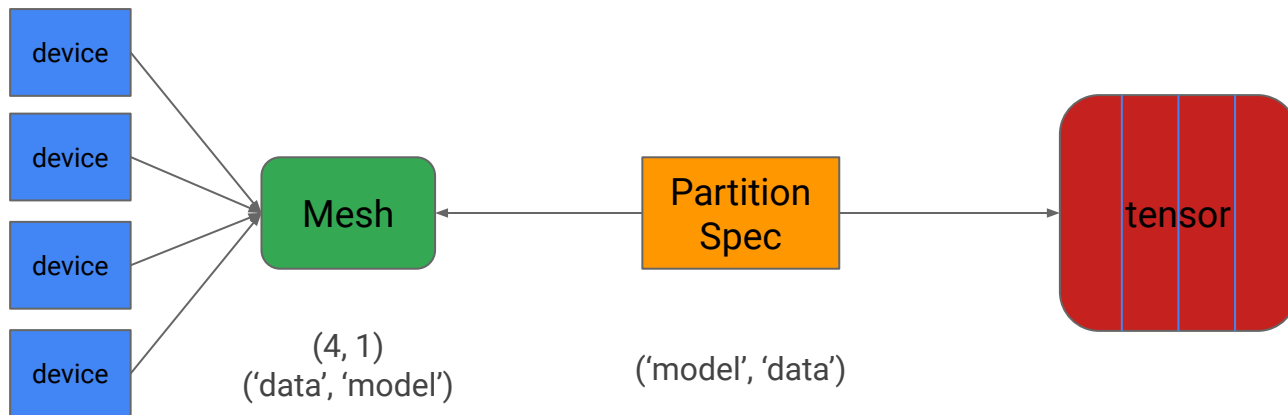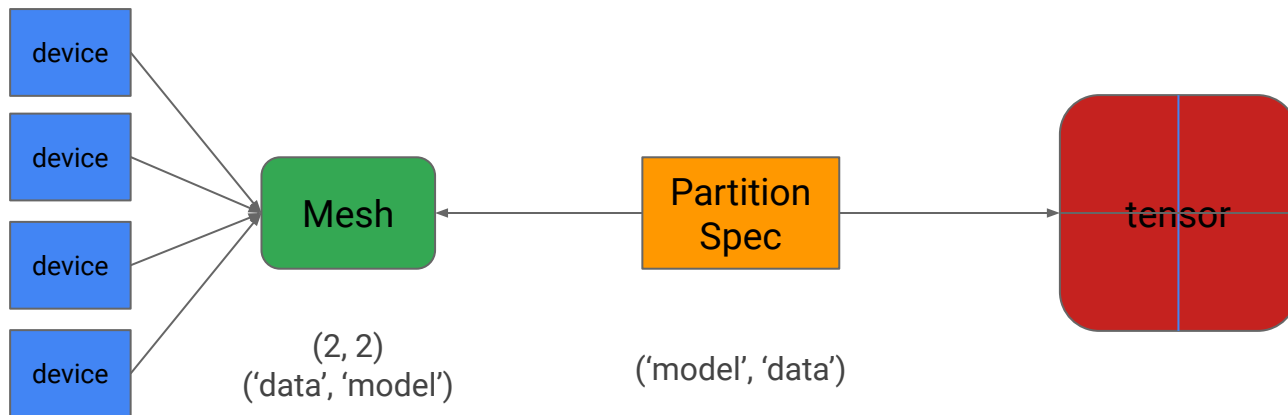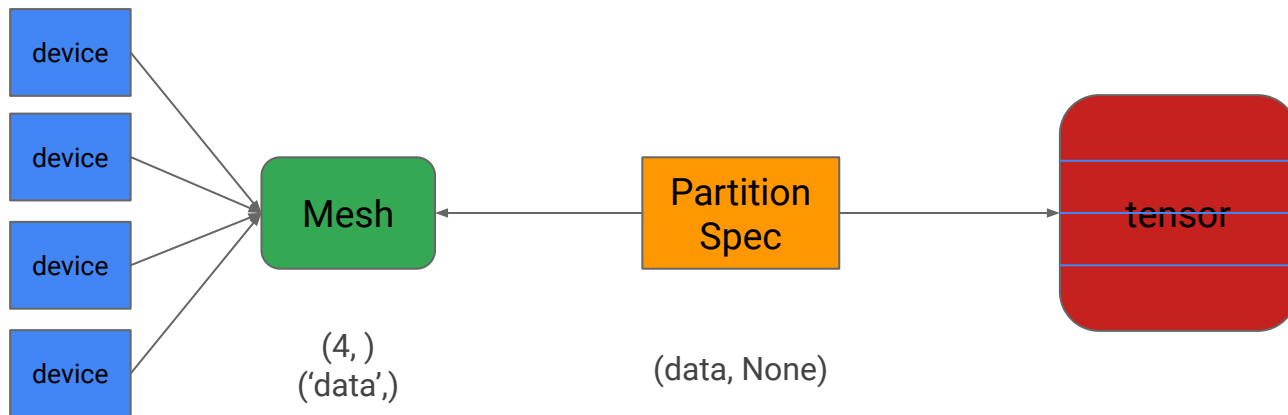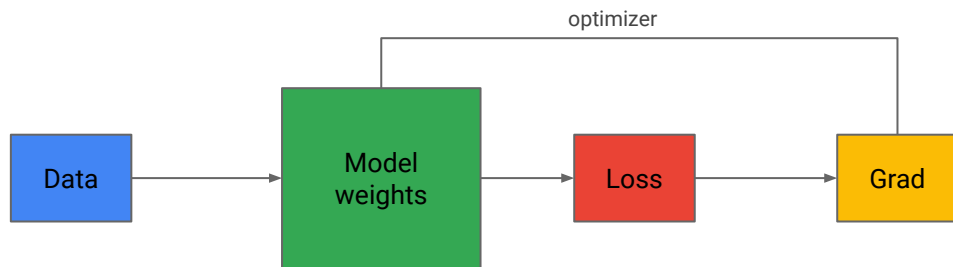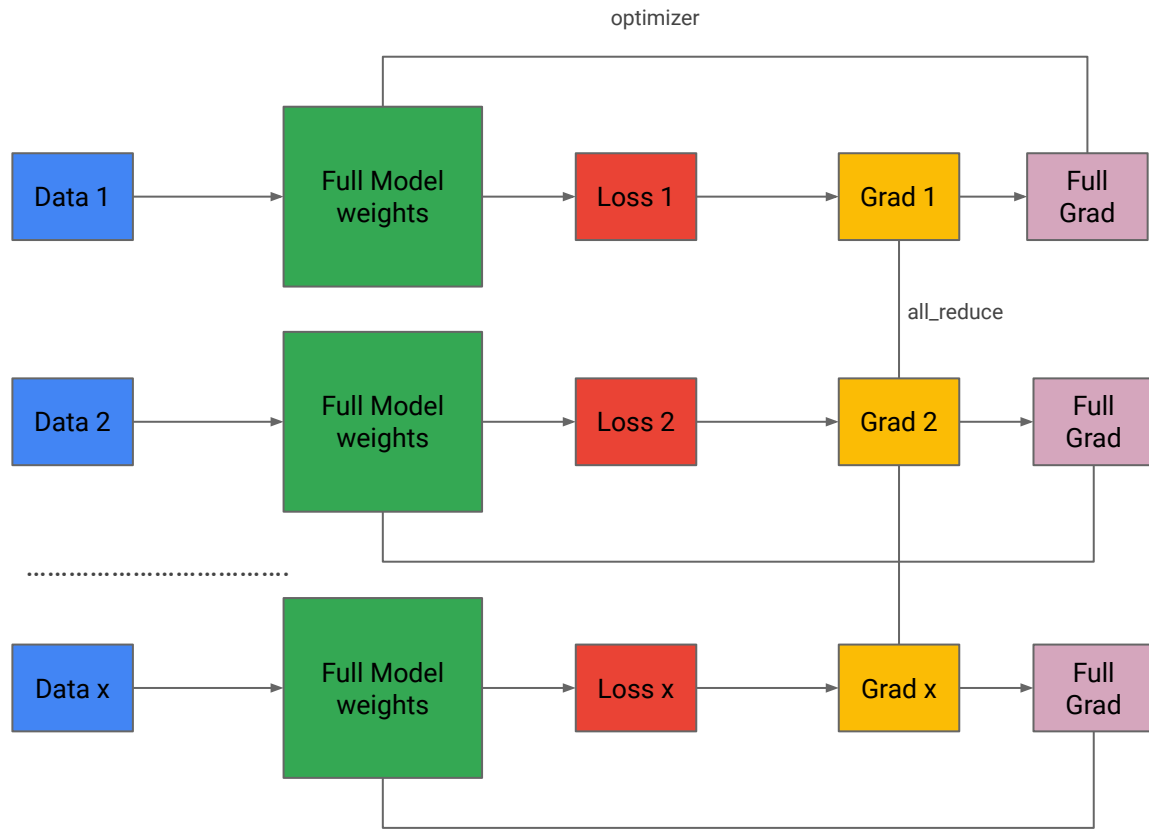
# Basic ideas

# Basic ideas

# Basic ideas

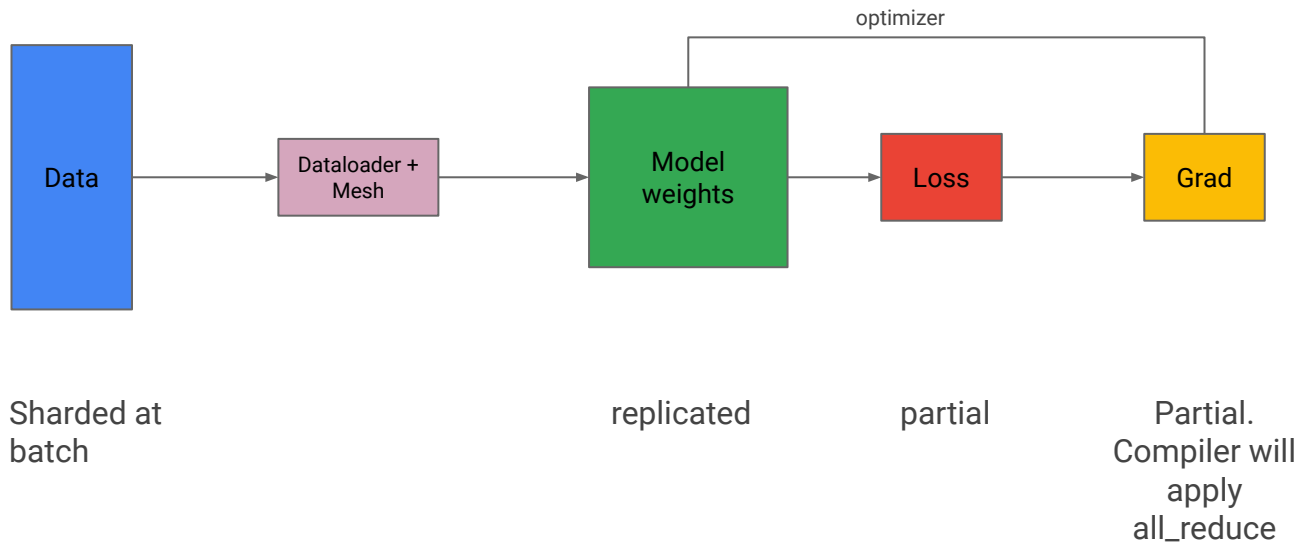# Basic ideas

# Basic ideas

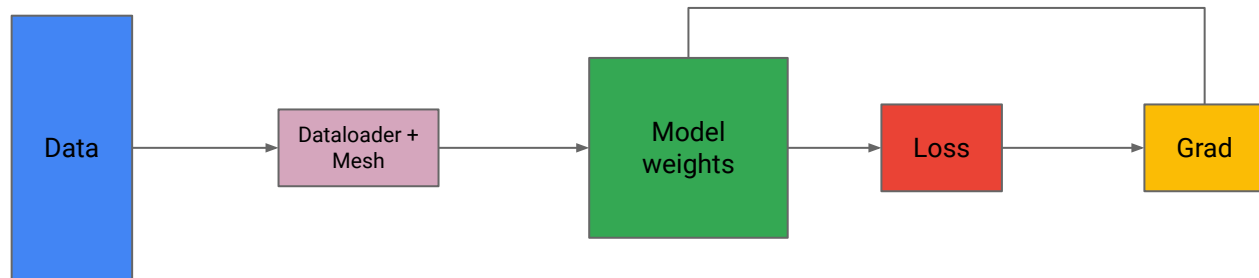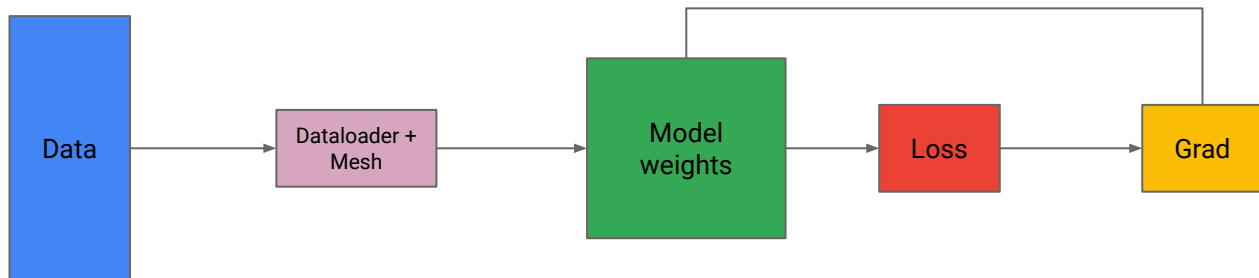# Single device

# Data Parallel

# SPMD + Data Parallel

# SPMD + Data Parallel

# SPMD + Data Parallel + multi host

# SPMD + Data Parallel + multi host