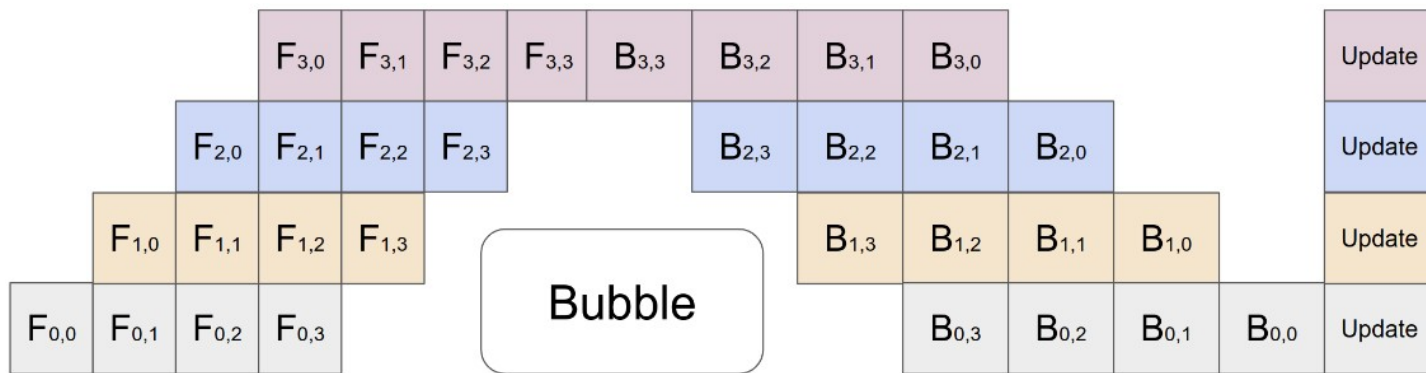
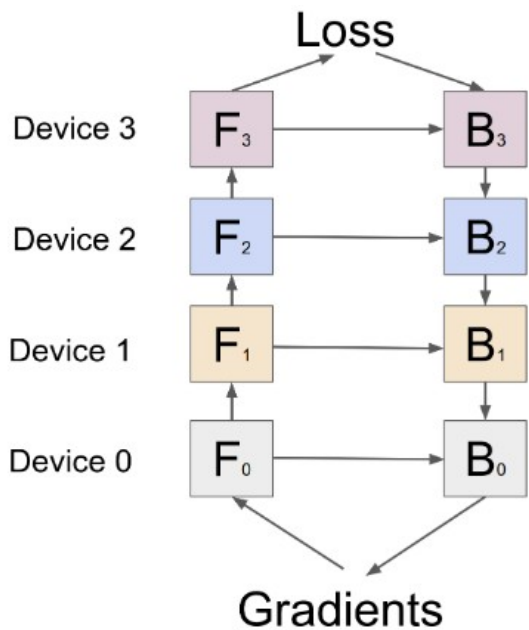


Pipelining

GPipe: arxiv.org/abs/1811.06965 – good starting point, *not* the 1st paper

Idea: split data into micro-batches and form a pipeline (right)



model size: $O(n)$
throughput: $O(n)$ – with caveats