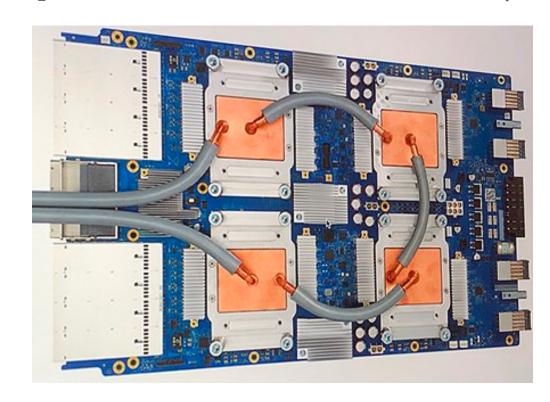
ARKADII

ONE SIDE TALK

About GNN in TPU and so on...

THIS IS VERY IMPORTANT FOR ALL MANKIND

- TPU is ASIC (high volume of lpc), higher ops/joule, higher bandwidth
- Al is a graph
- Graphs are large (no way to use GNN directly)
- Not too large
- Hardware optimal execution is high profitable (Google used autotuner for XLA, a machine learning compiler, discovered 10–20% speedup on SOTA models)
- Autotuner is expensive
- Compilers need scoring



[2,4,16]

{1,0,2} \(\sqrt{} \)

conv [2,8]

{0,1} **{1,0**}

{1,0,2}

[128]

[4,16,8]

{0,2,1}

{2,0,1}

