

# ***Square: Towards Sharing Internet-Scale GPU Clouds Fairly and Efficiently***

**Ningxin Su, Freeman Cheng, Baochun Li**

**Ningxin Su, Assistant Professor  
Information Hub, IoT Thrust  
HKUST (GZ)**

**IWQoS 2026, Istanbul**

## Google DeepMind

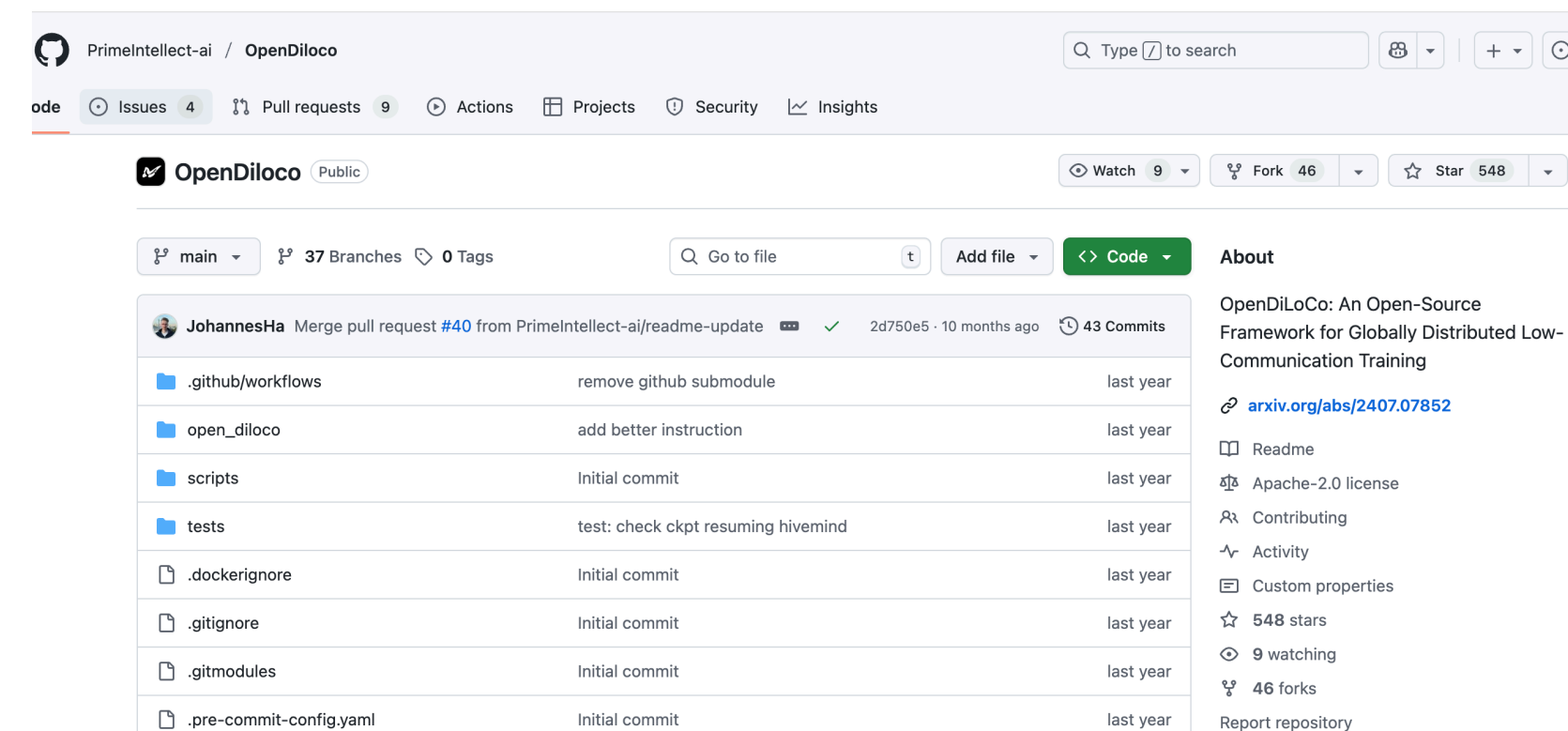
- DiLoCo
- a variant of FedAvg



<https://arxiv.org/abs/2311.08105>

## PrimeIntellect-ai

- OpenDiLoCo
- Open-Source

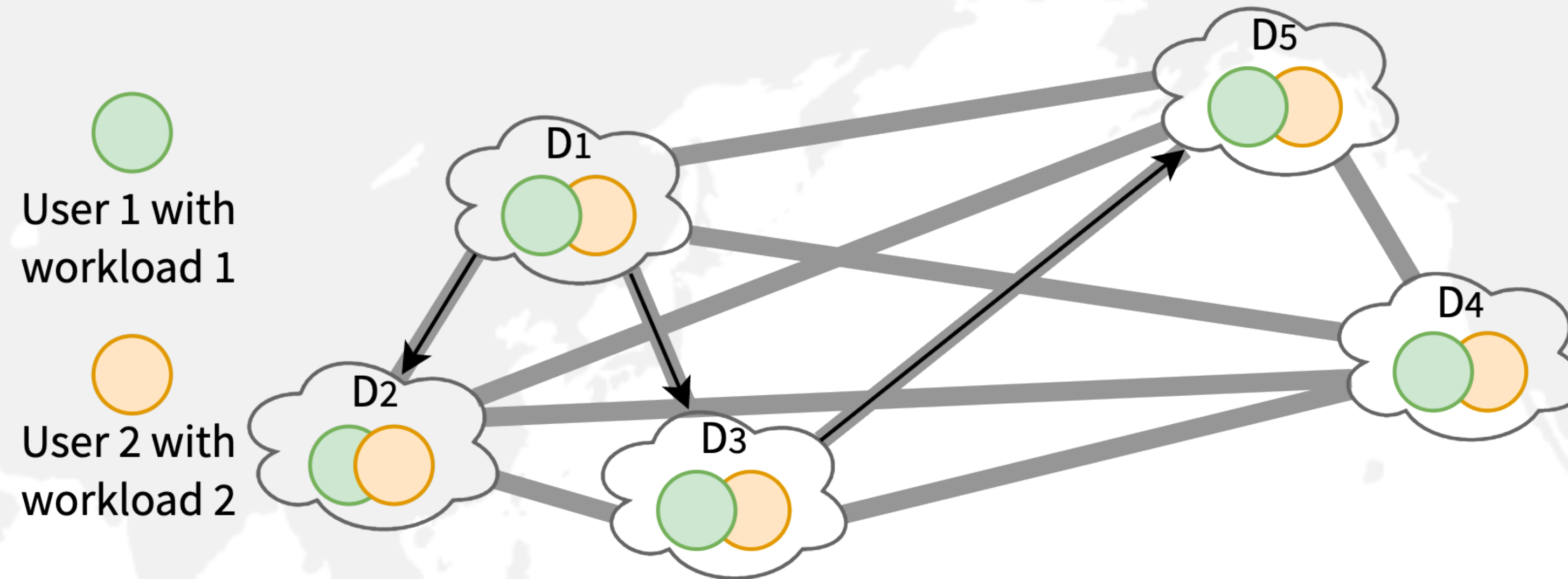


<https://github.com/PrimeIntellect-ai/OpenDiloco>

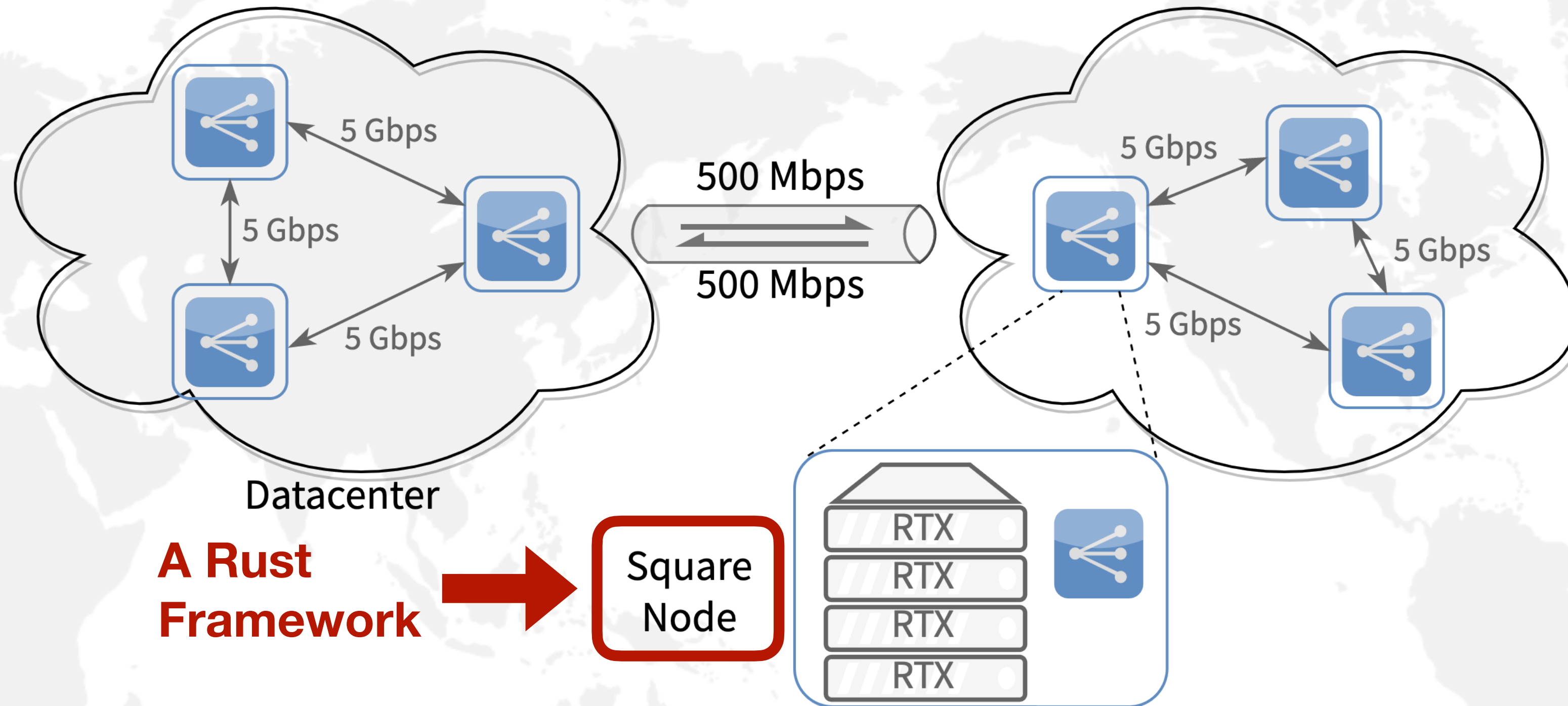
## Gongji Cloud



<https://suanli.cn/>

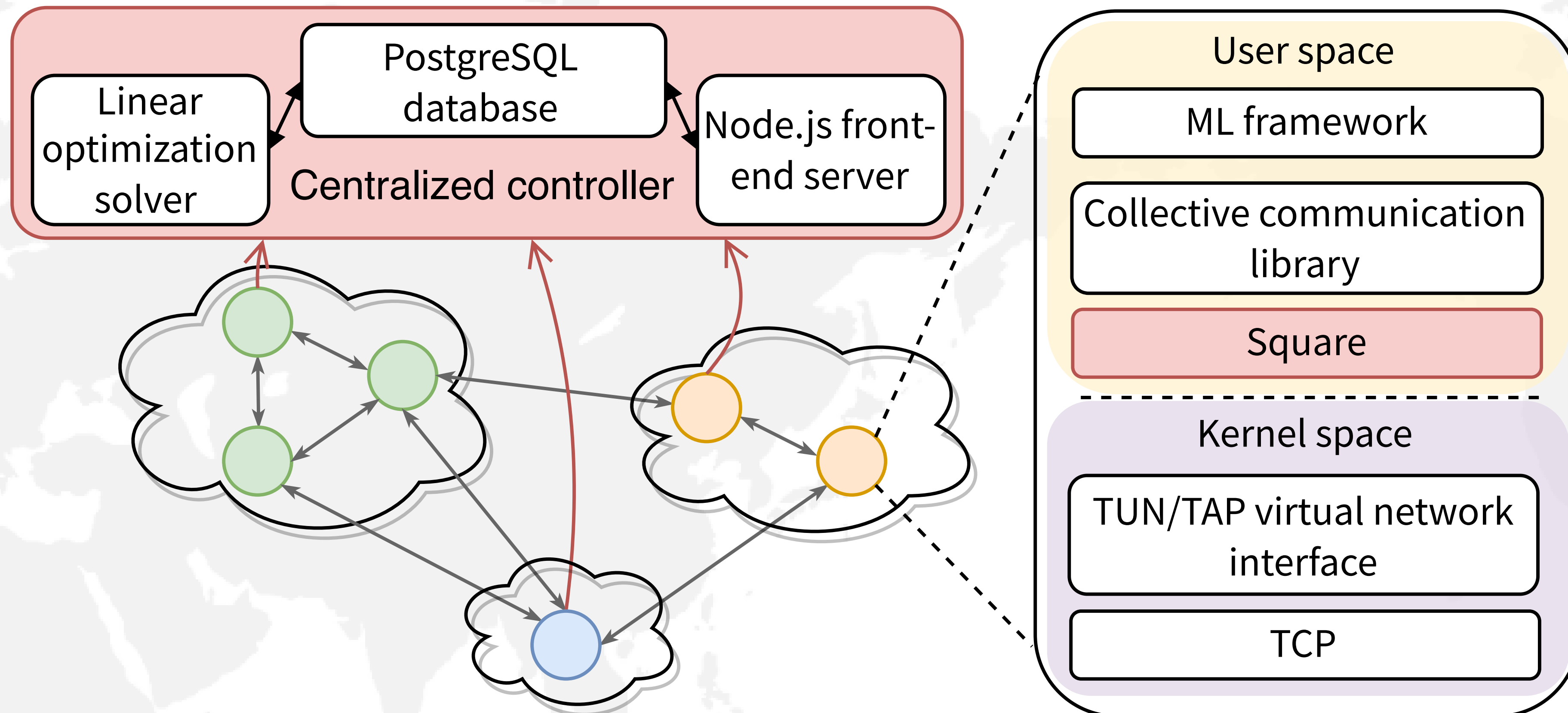


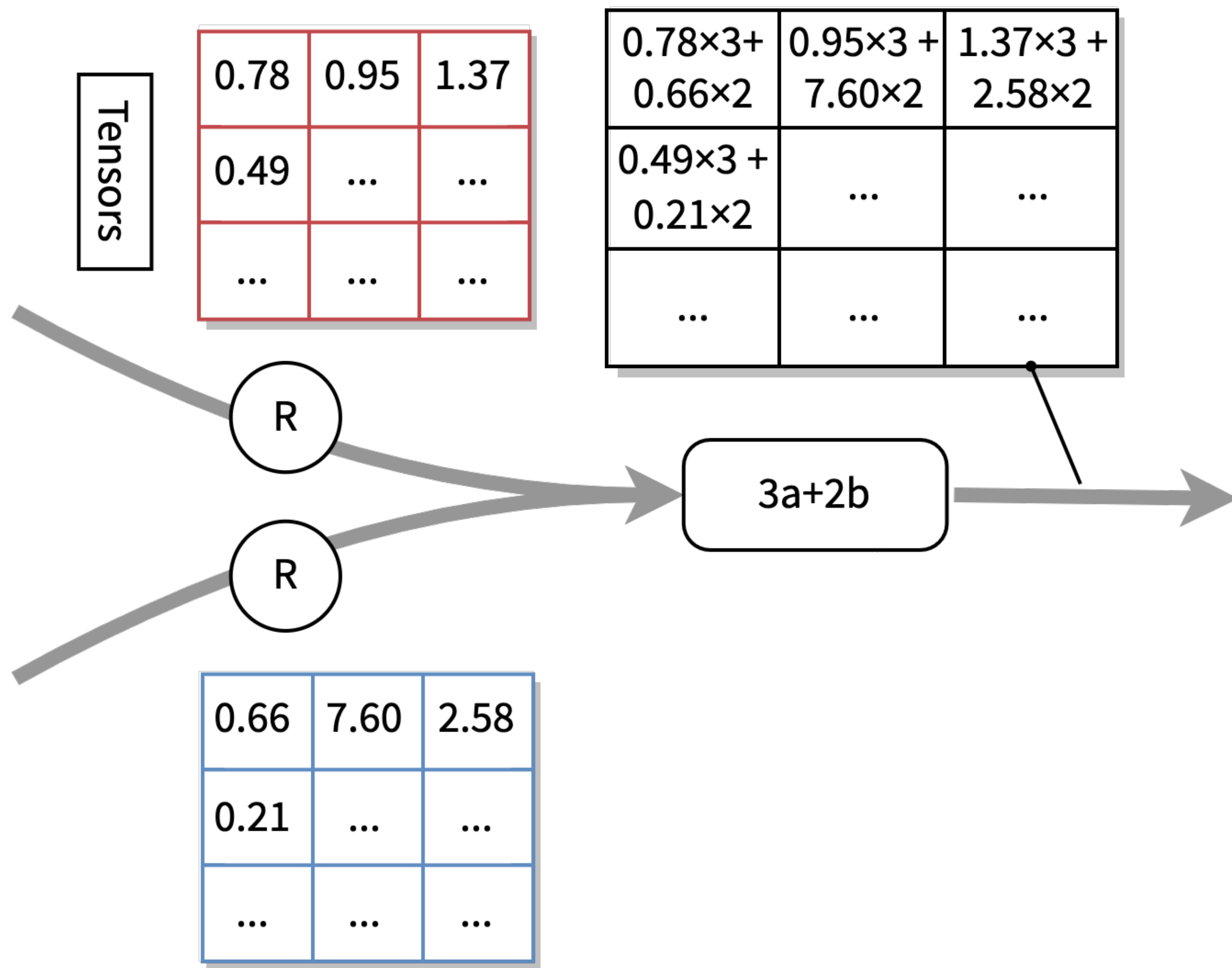
**When many users train across many cloud regions, how should scarce WAN bandwidth be shared fairly and efficiently?**

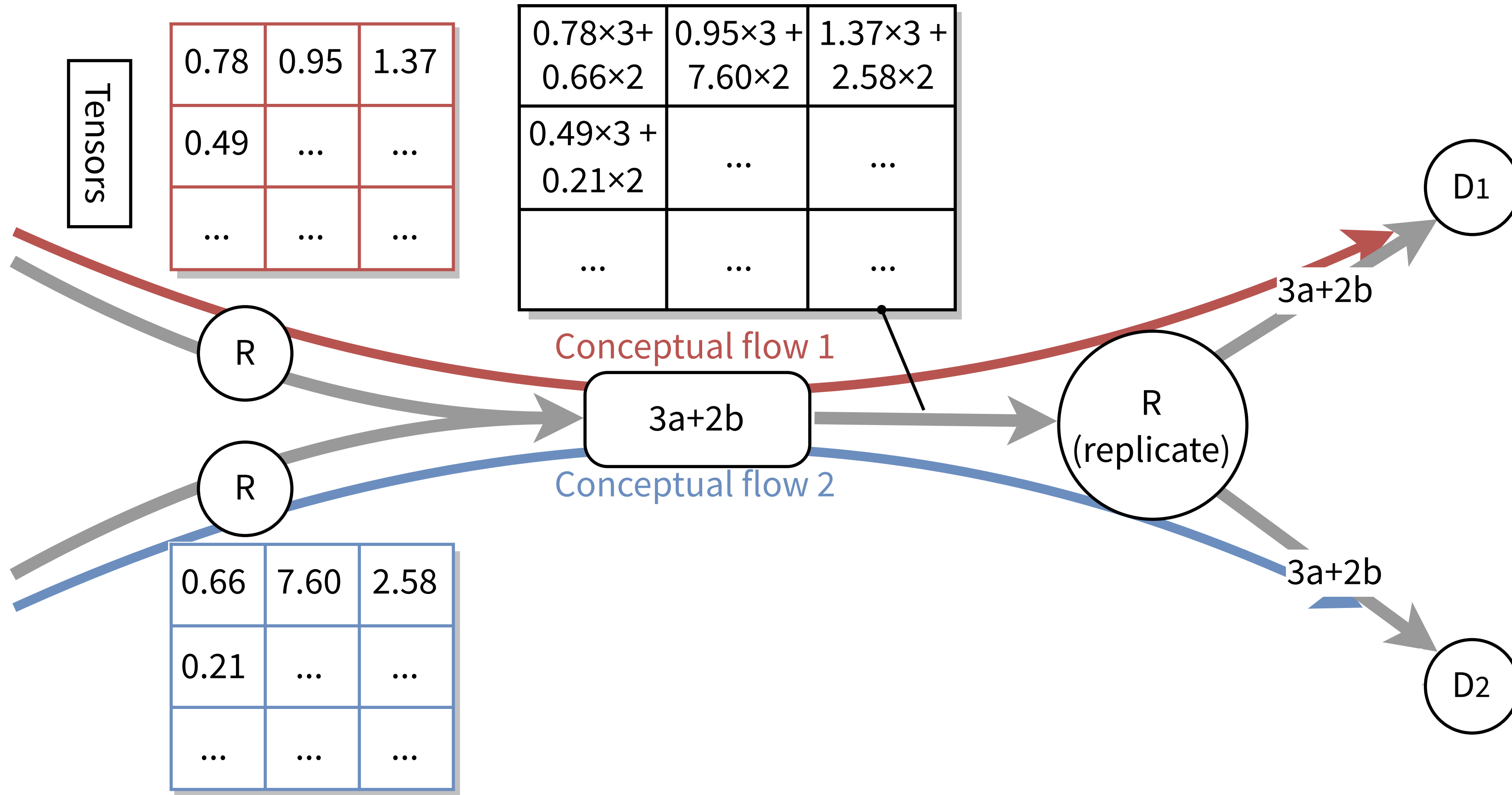


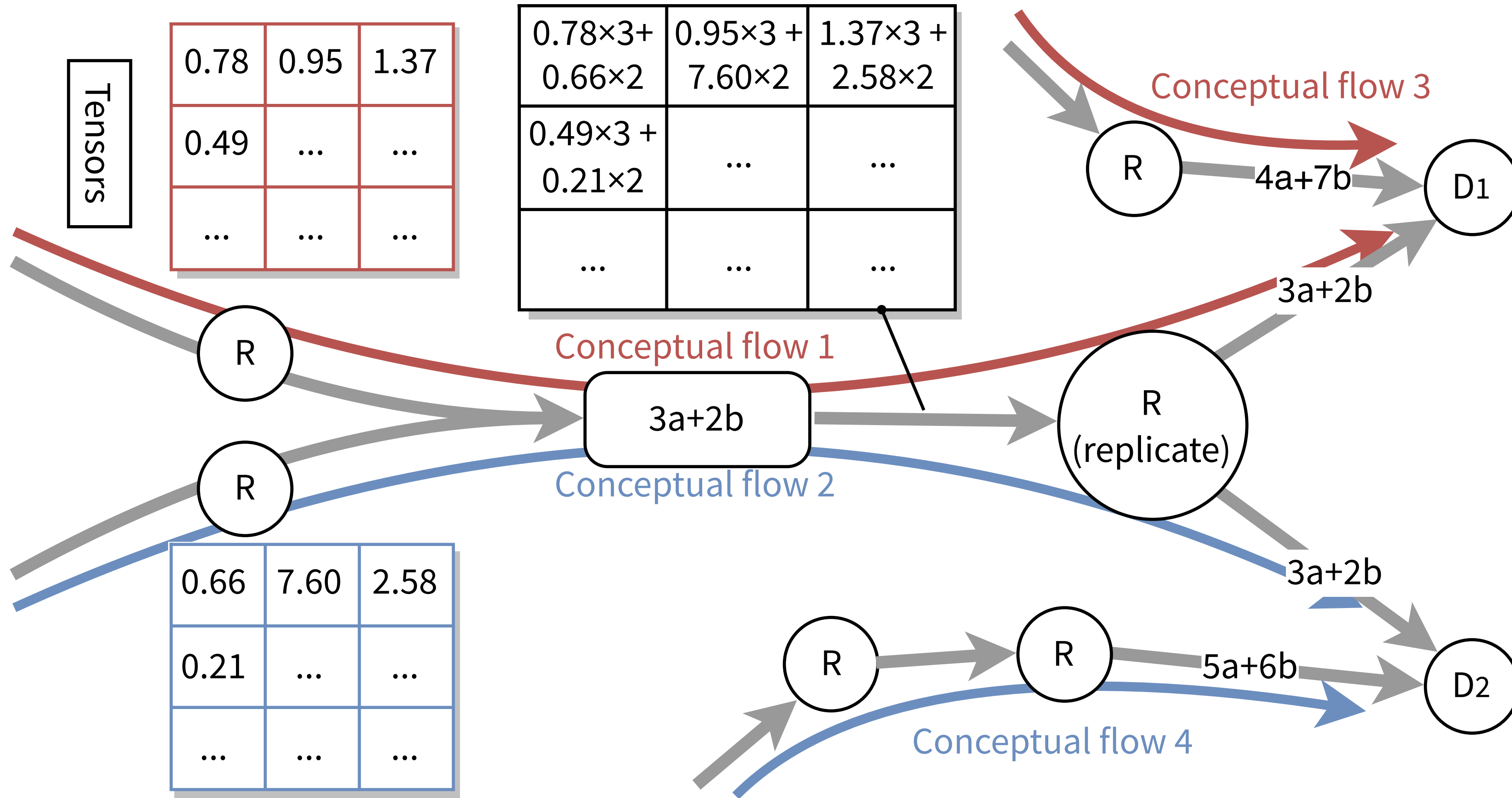
**A Rust  
Framework**

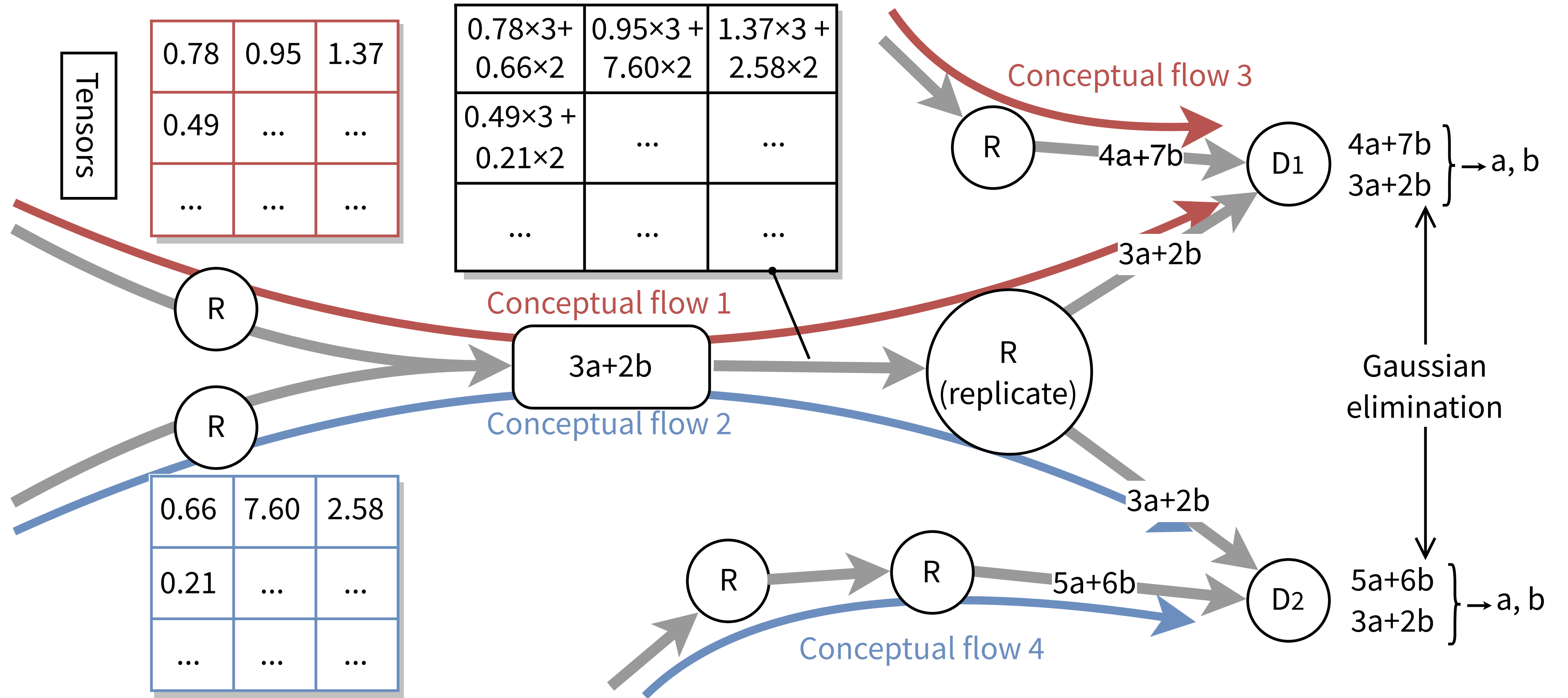
**Square  
Node**

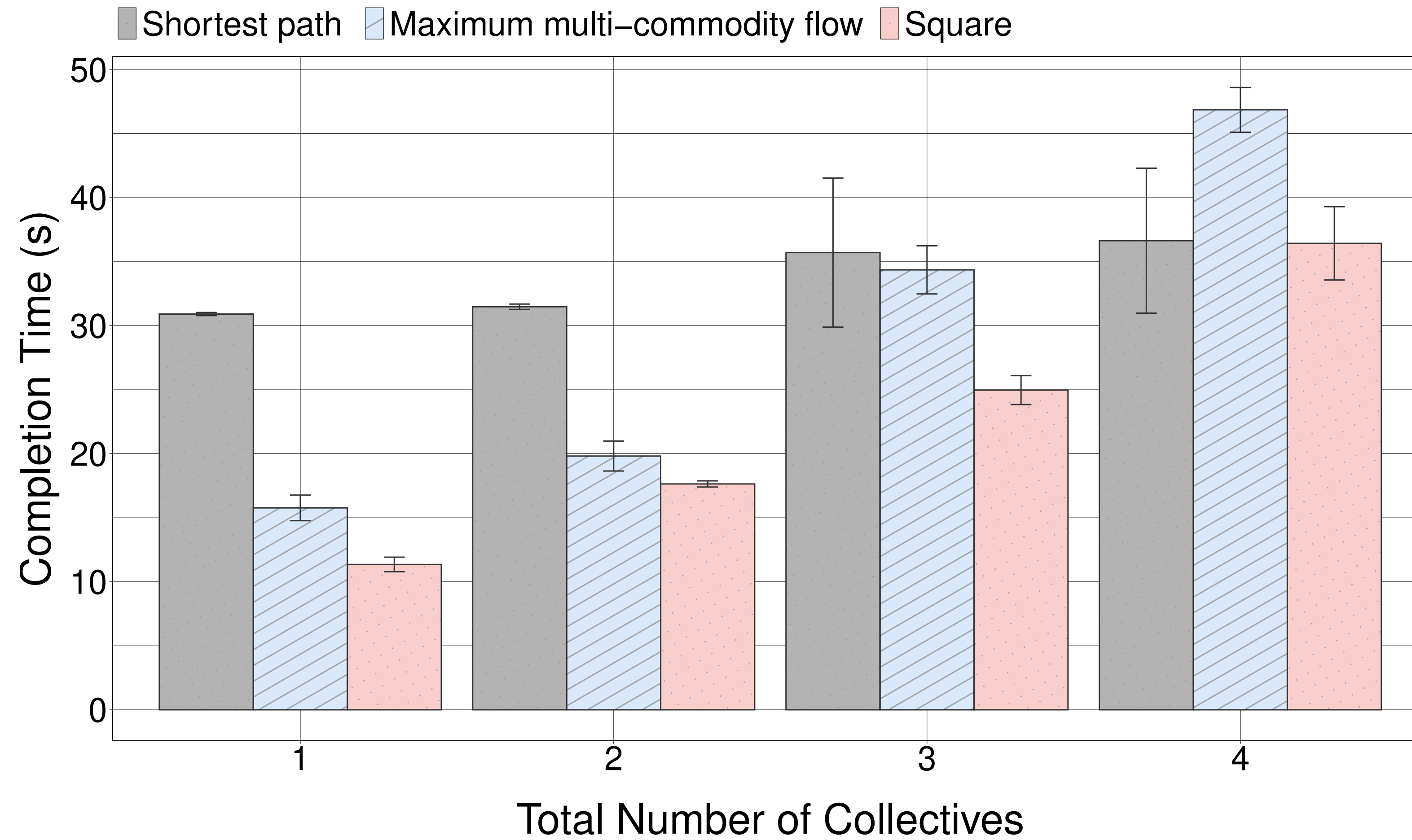












Experimental results over a real Digital Ocean inter-datacenter overlay.

# Thank you!

[ningxinsu@hkust-gz.edu.cn](mailto:ningxinsu@hkust-gz.edu.cn)

[nebulis-lab.com](http://nebulis-lab.com)

[ningxinsu.github.io](https://ningxinsu.github.io)