MX+: Pushing the Limits of Microscaling Formats for Efficient Large Language Model Serving

Jungi Lee, Junyong Park, Soohyun Cha, Jaehoon Cho, Jaewoong Sim Seoul National University



