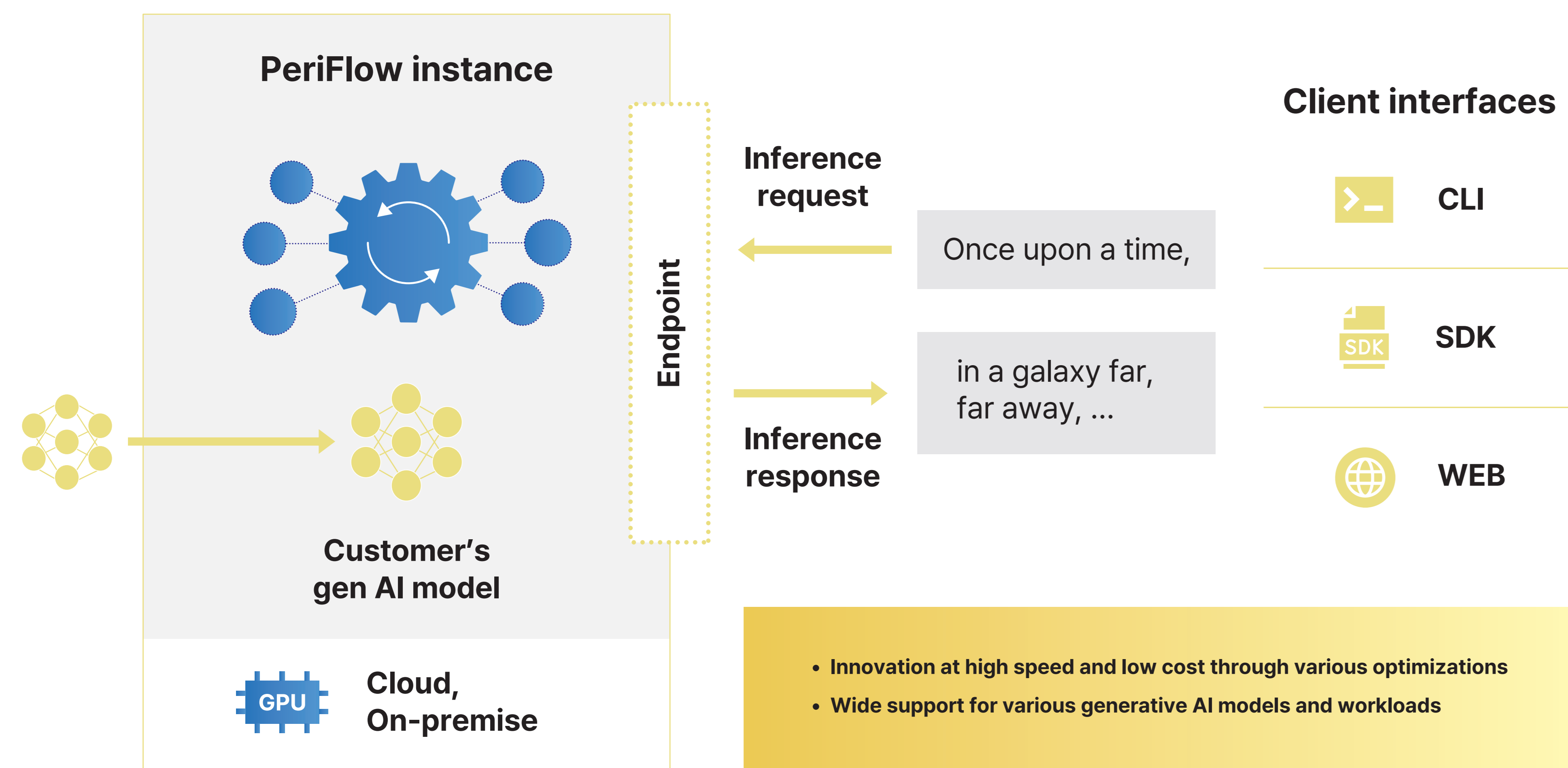
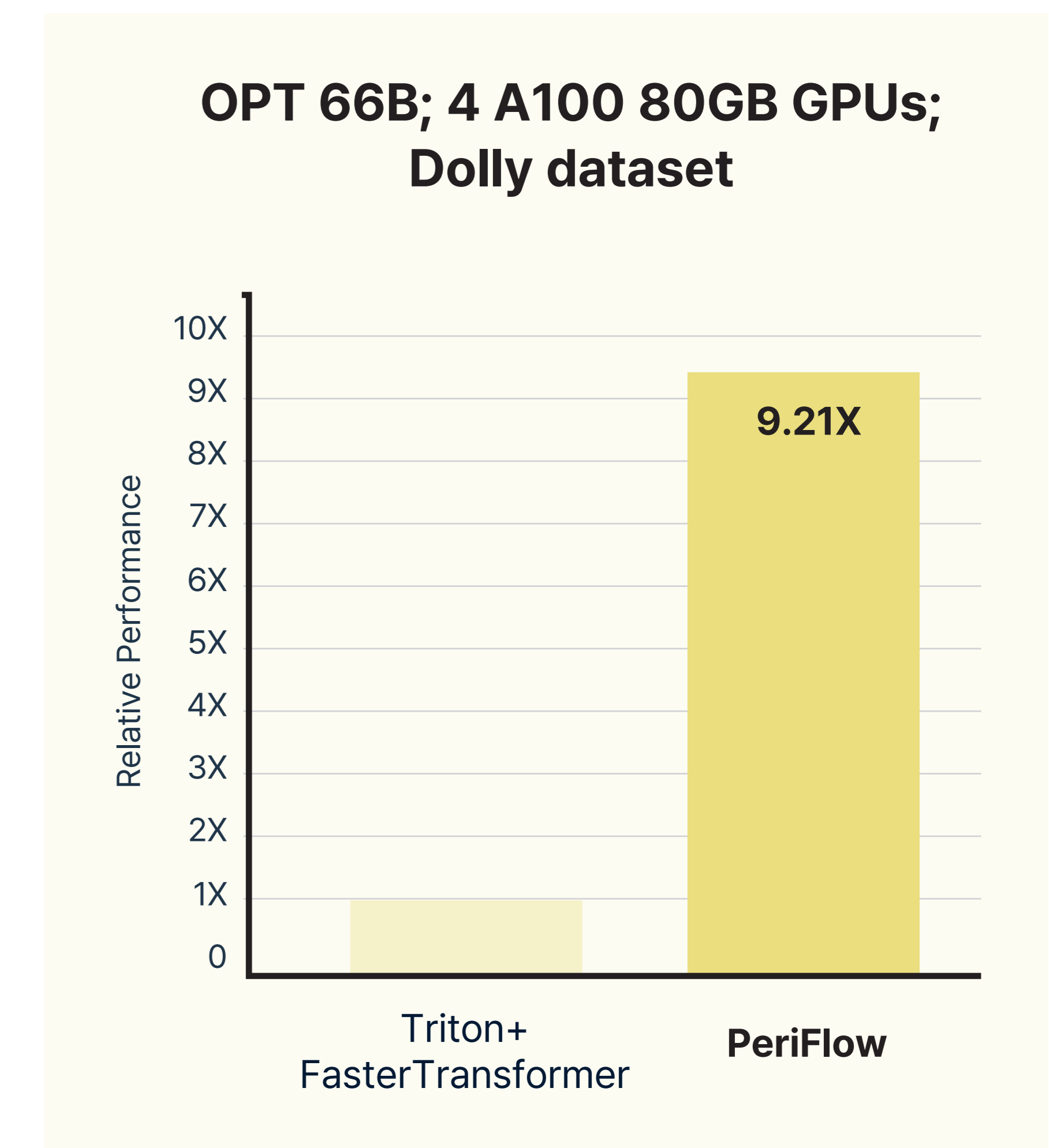


PeriFlow Engine for serving generative AI models

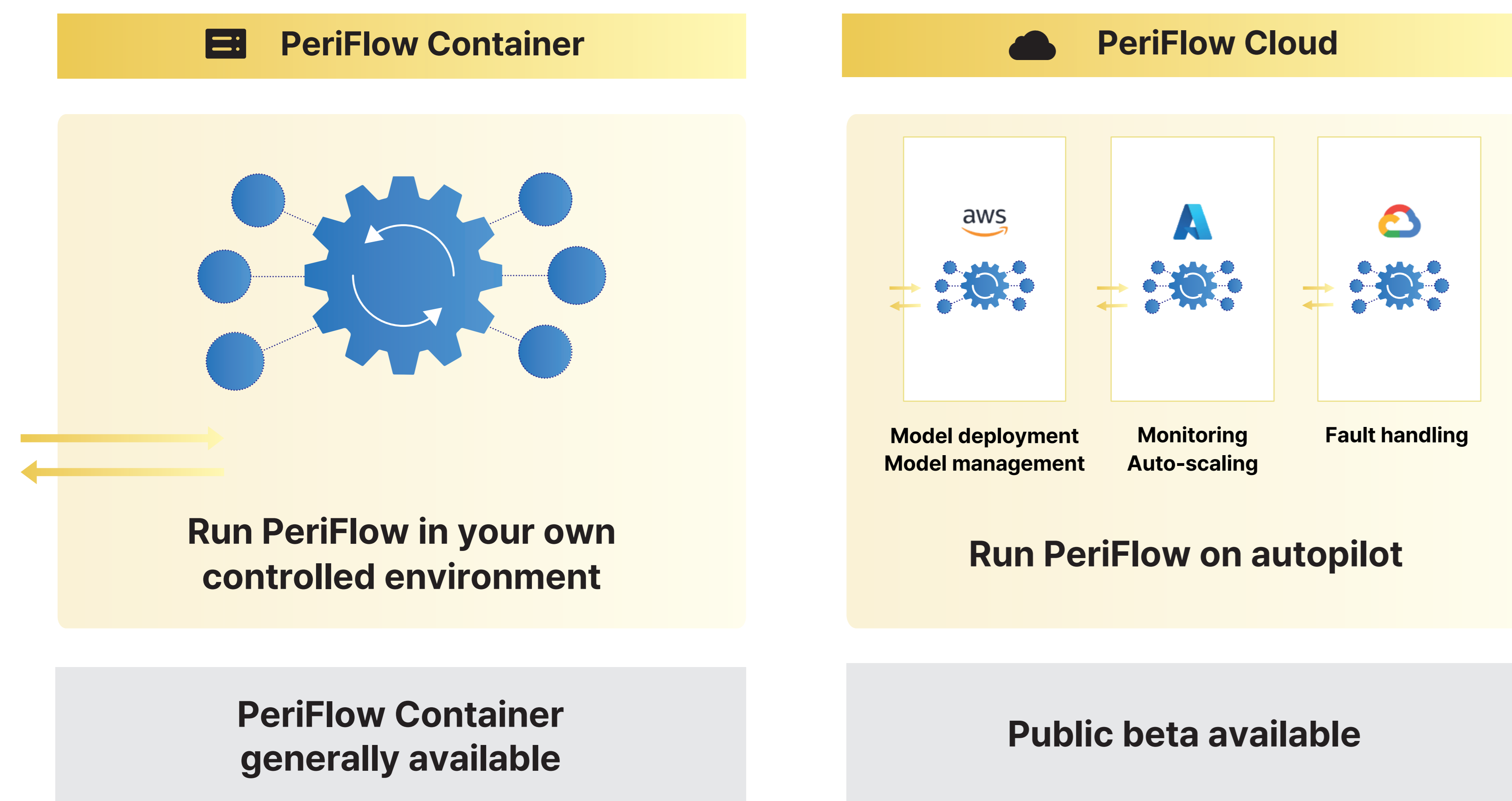


PeriFlow Engine for serving generative AI models

- 10x faster than Triton and FasterTransformer
- Reduce GPUs needed by 10x → reduce costs, improve reliability
- Wide support for various gen AI models and workloads
- Streamline gen AI inference optimization and deployment



Two ways to use PeriFlow : PeriFlow Container and PeriFlow Cloud



PeriFlow Cloud interface

