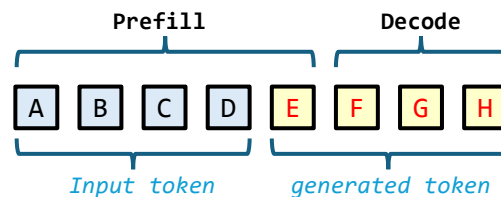
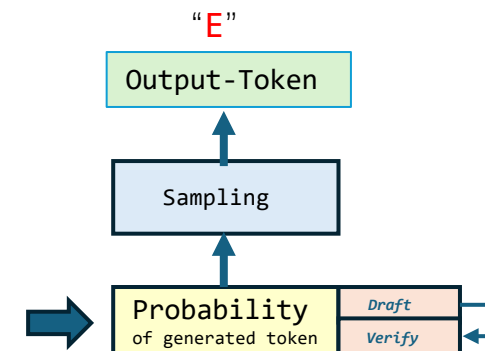
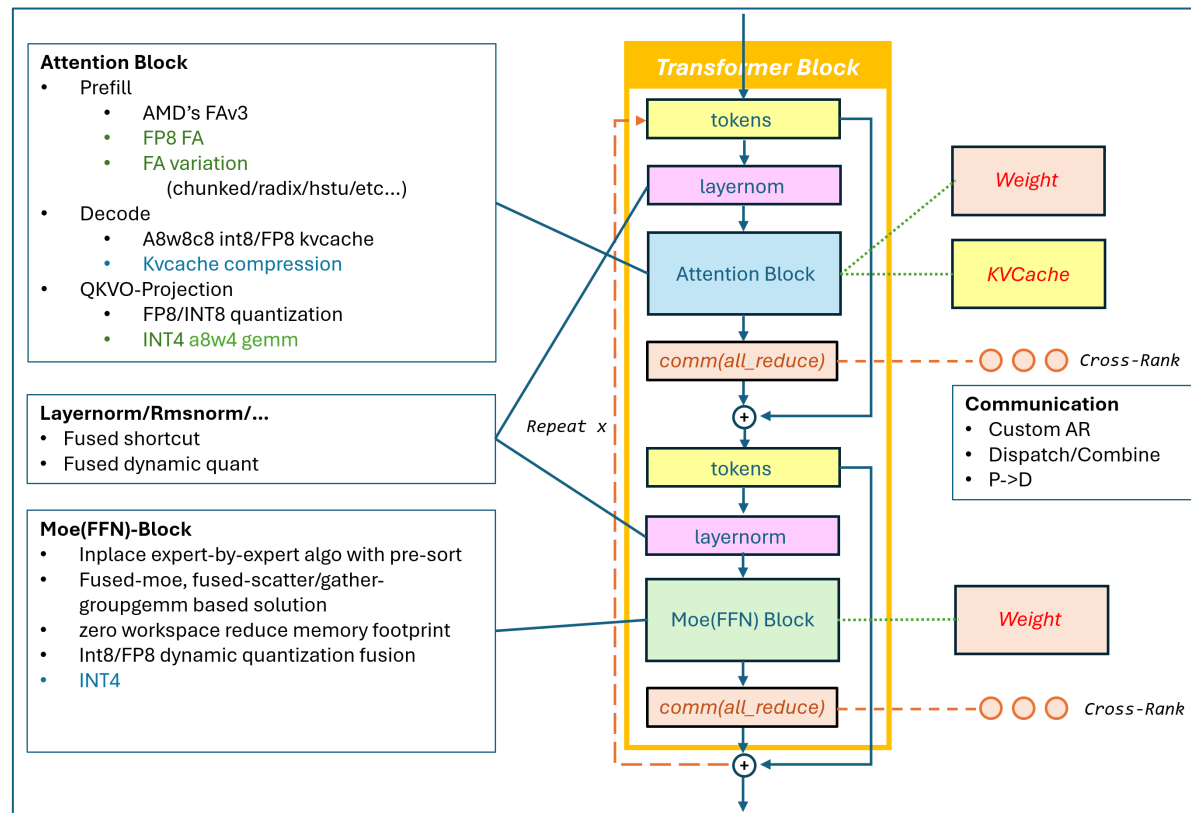
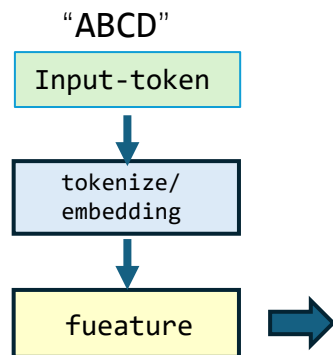


AITER / MoRI

Introduction



carlus.huang@amd.com

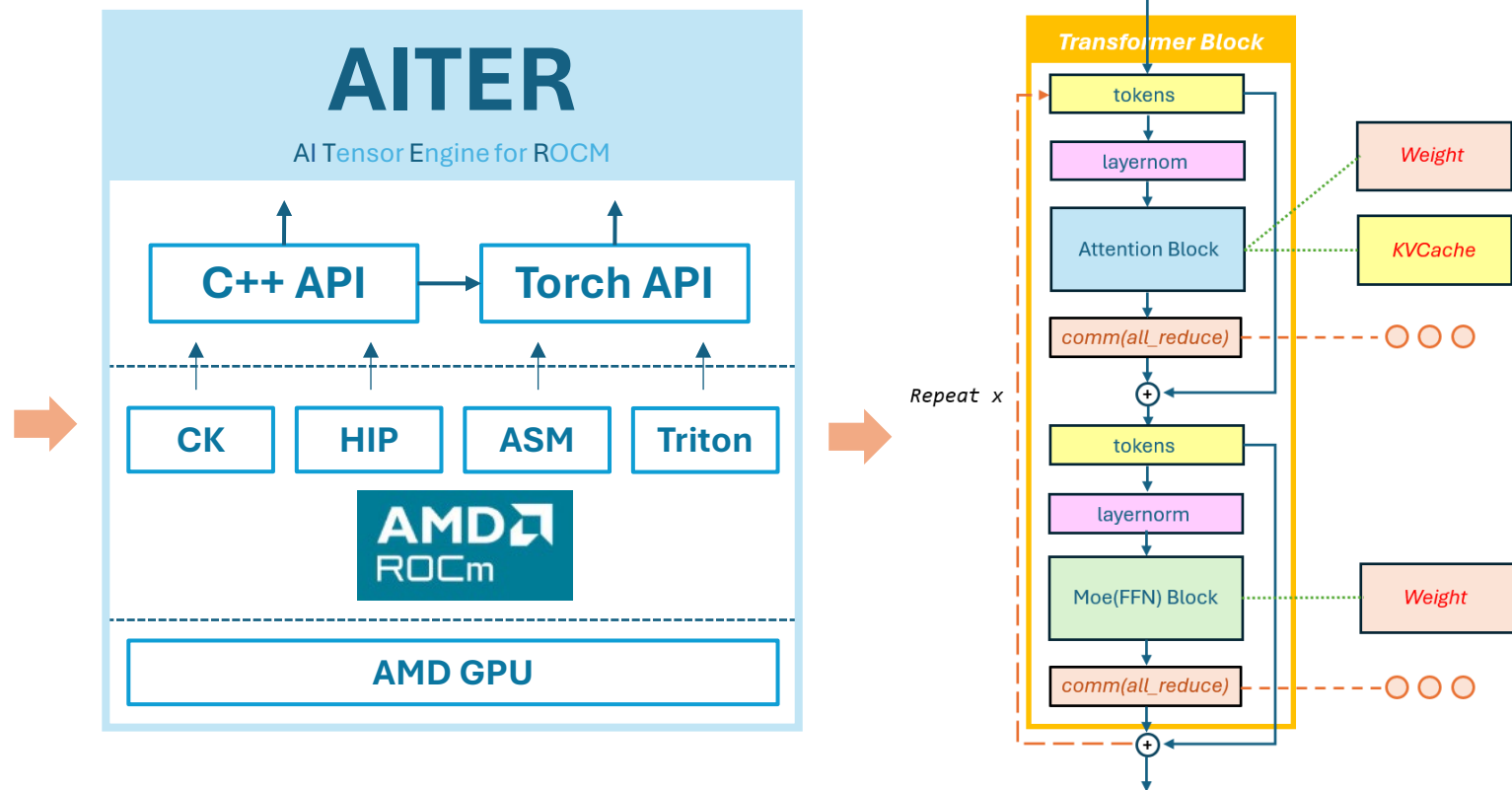


=> View LLM as Building Block

AITER overview

GITHUB : ROCm/aiter

Prefill Attention	F/B	FAv3 FWD FP16/BF16 FAv3 BWD FP16/BF16 MLA FP16/BF16 FA FP8 FWD+BWD (Block-Scale) Chunked-prefill
Decode Attention	F	Paged Attention FP16/BF16 Paged Attention FP8 per-tensor quant Paged Attention FP8/INT8 with KV per-token quant KVCache update & Rotary Batched Decoding MLA Decoding
Fused-Moe	F	Moe-Sorting kernel and tiling solution FP16/BF16 per-token Fused-Moe FP8/INT8 per-token Fused-moe FP8 per-tensor Fused-Moe FP8/INT4 per-tensor Fused-Moe Fused-FFN
Low Precision Gemm	F	FP8 per-token/channel Gemm FP8 Block Scale Gemm INT4 weight only gemm
Distributed Gemm	F/B	(need evaluation)
Normalization and fusion	F	layernorm+quant/shortcut rmsnorm+quant/shortcut
Custom Comm.	F/B	AR/AG... fused with noamalizatoin AR/AG... quantized Optimized hipgraph support
Conv2d/2d	F/B	fp16/bf16 fwd/bwd/wrw fusion with bias/activation, etc..
And more fused ops...		



AITER scope:

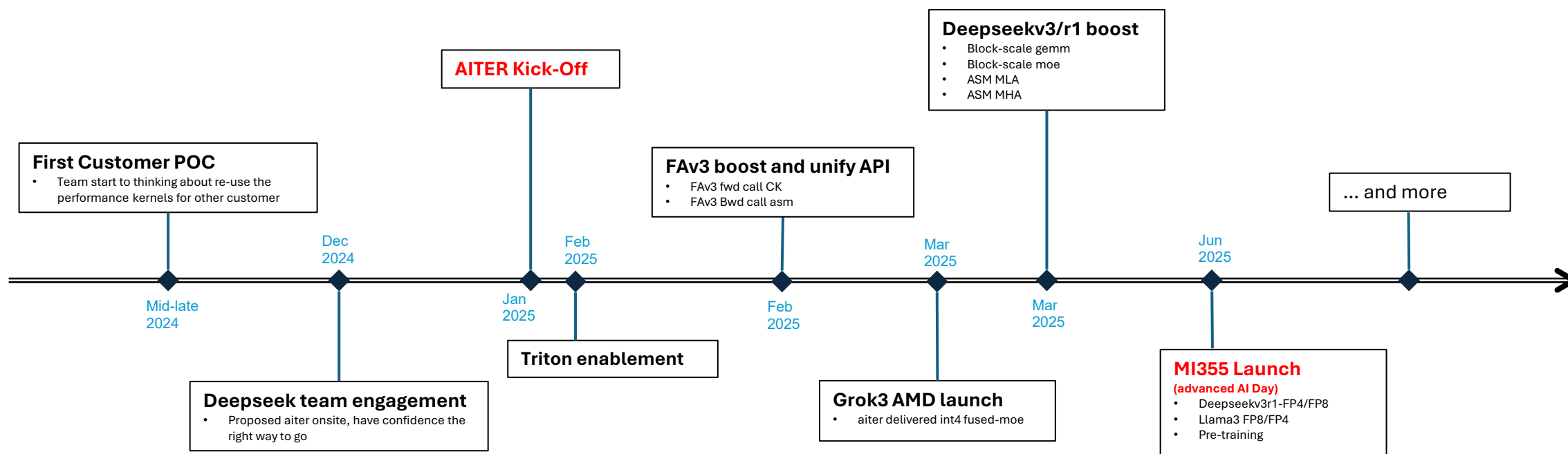
- A lightweight and customizable/tunable/deployable AI tensor engine.
- A collection of inference AI operators
- A collection of training AI operators
- Fused dynamic quantization kernels and offline utilities.
- LLM block-level solutions to plug into customer's framework. (fused-moe, quant-KVCache, etc...)
- C++/Pytorch API and package
- heuristic/user-driver-tuning for self-hosted kernels written in CK/ASM/HIP/Triton

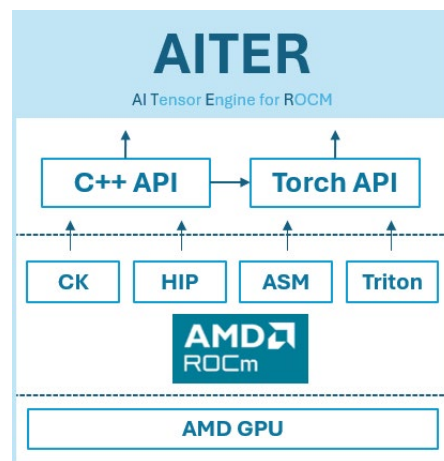
AITER is not:

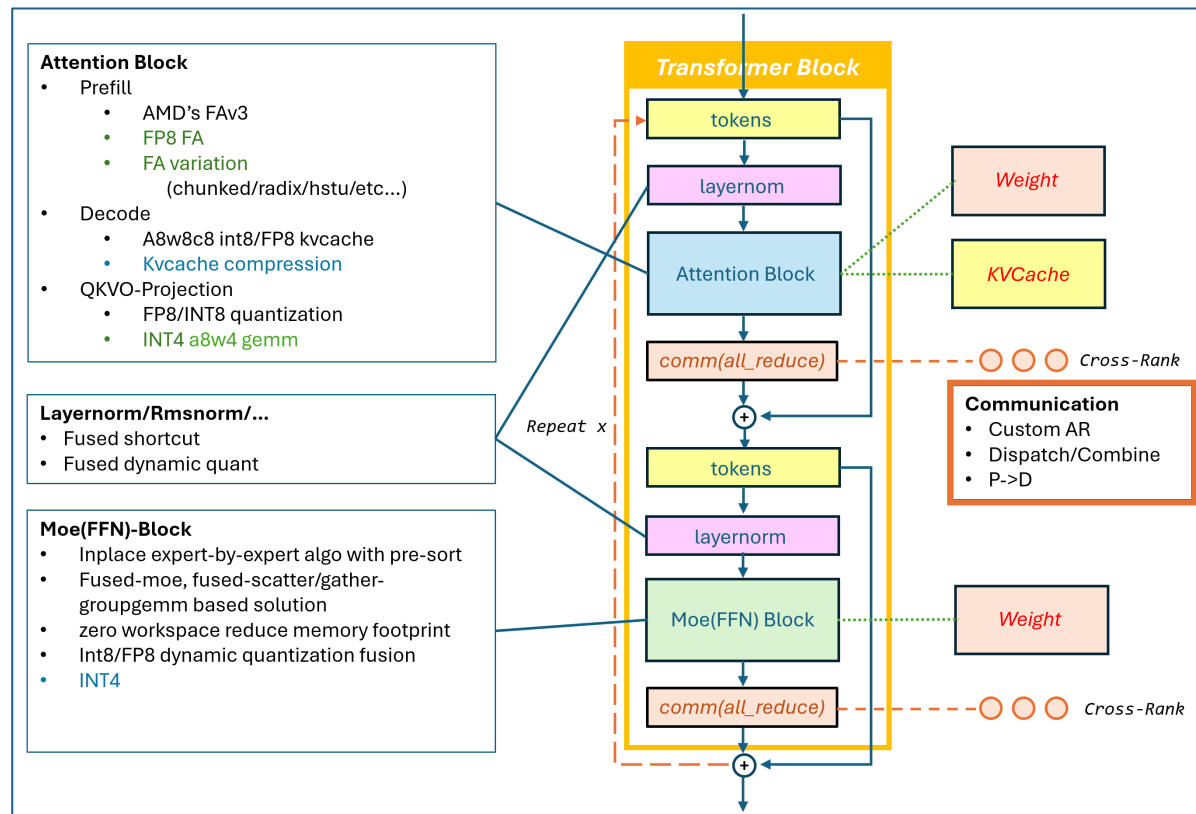
- a LLM service engine or LLM workload balance scheduler
- a distributed inference or training framework
- a graph compiler
- a thin wrapper on top of existing rocm libraries.

ALTER timeline

- Hosting reusable optimized kernels
- Establishing standardized integration frameworks
- Enabling rapid deployment across clients





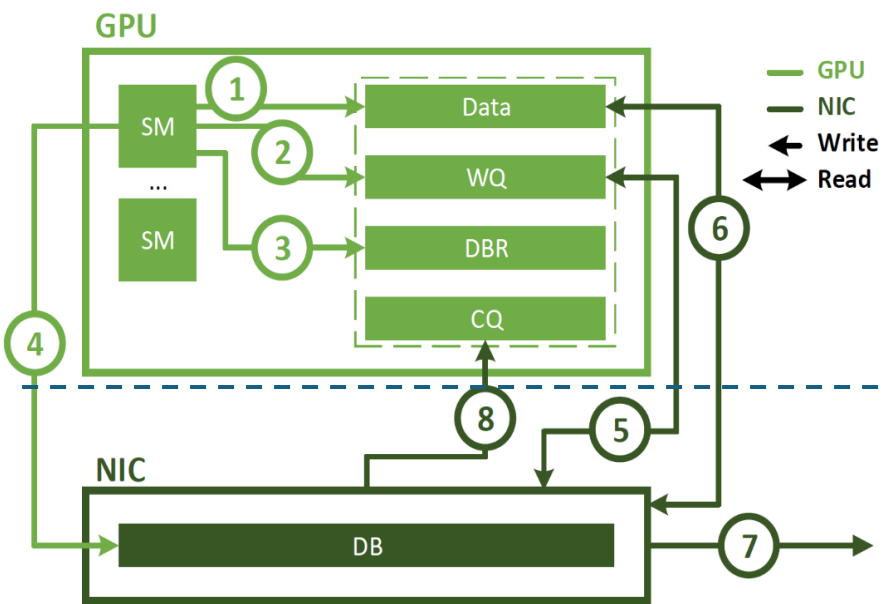


MoRI
 (Modular RDMA Interface)
 LLM Native Communication Engine

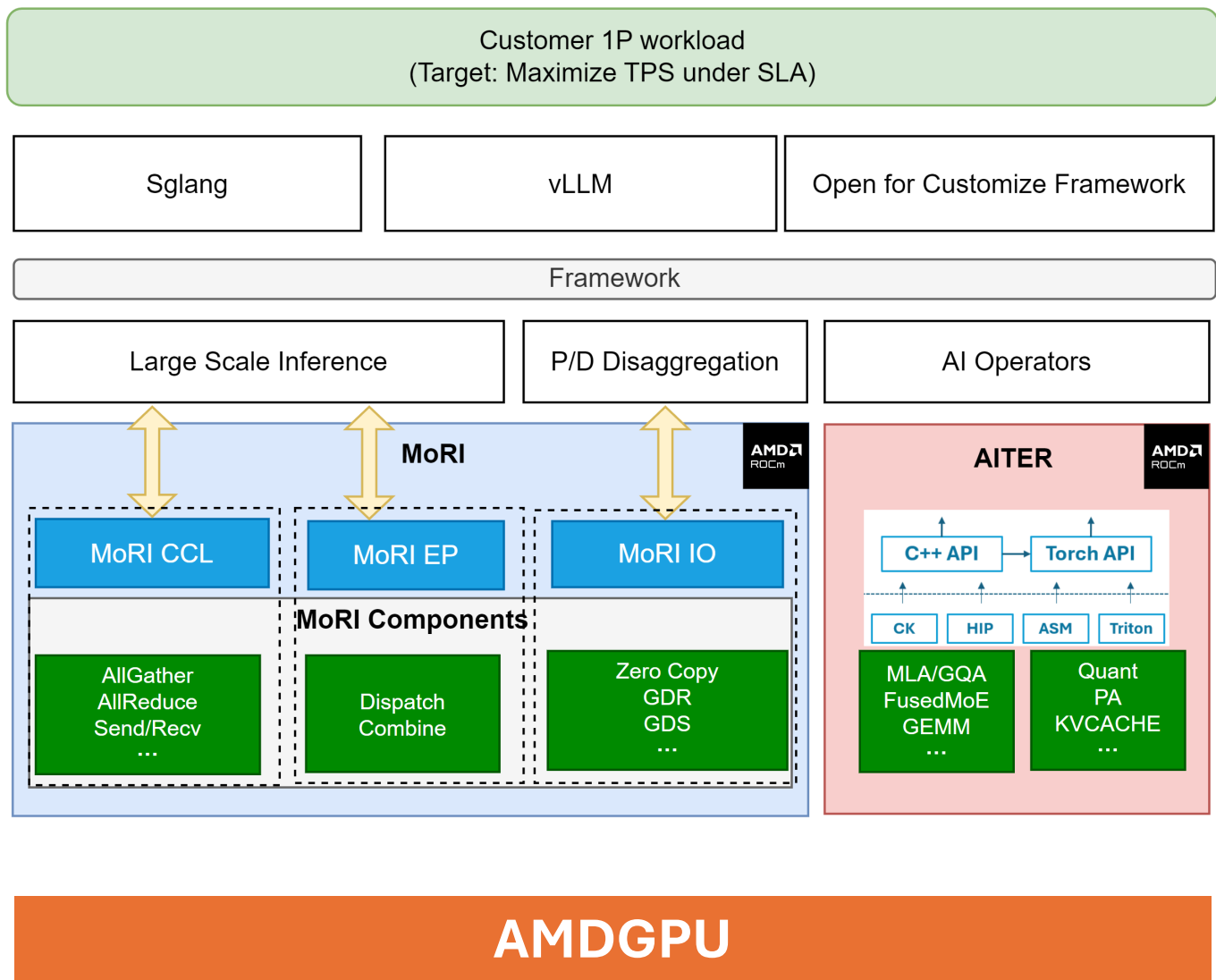
GITHUB : ROCm/mori

MoRI is all you need to implement high-performance communication on GPGPUs.

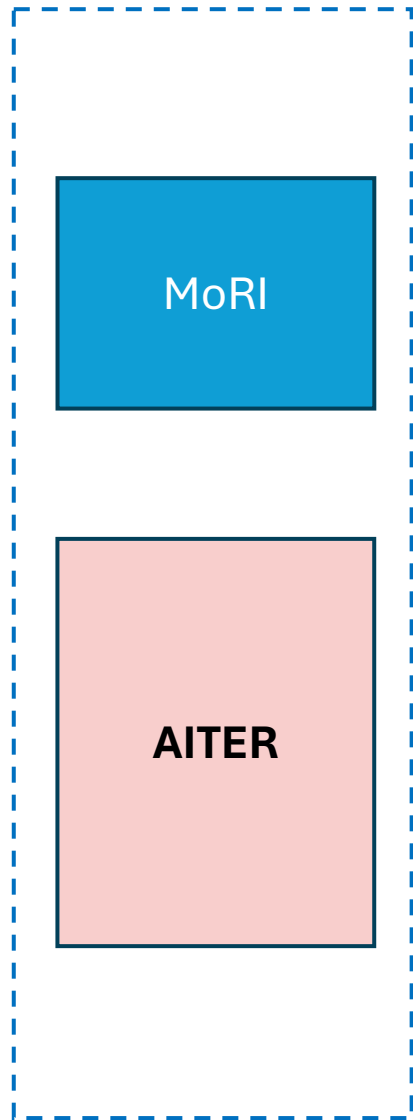
MoRI = DeepEP + shmem + CCL + ...



MoRI			
application	Custom Kernels (e.g DeepEp / pplx-kernel)	SHMEM API	
	CCL	Fused Kernels	
framework	Transport Context Management		Buffer Management
	Topology Detection		...
core	thread-level	wave-level	block-level
	data transfer	synchronization	signal & atomic reduction & quantization
transport	IBGDA	GPU P2P	NVMe
Hardware	Broadcom	Mellanox	Pensando

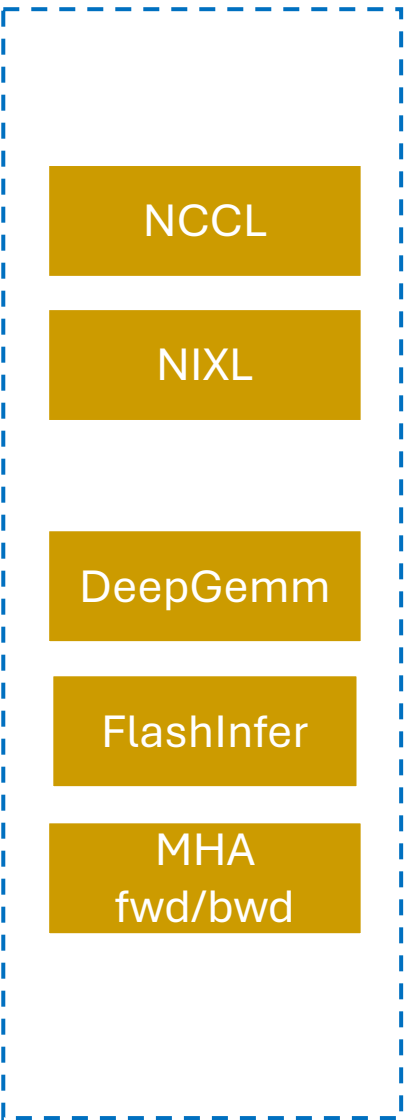


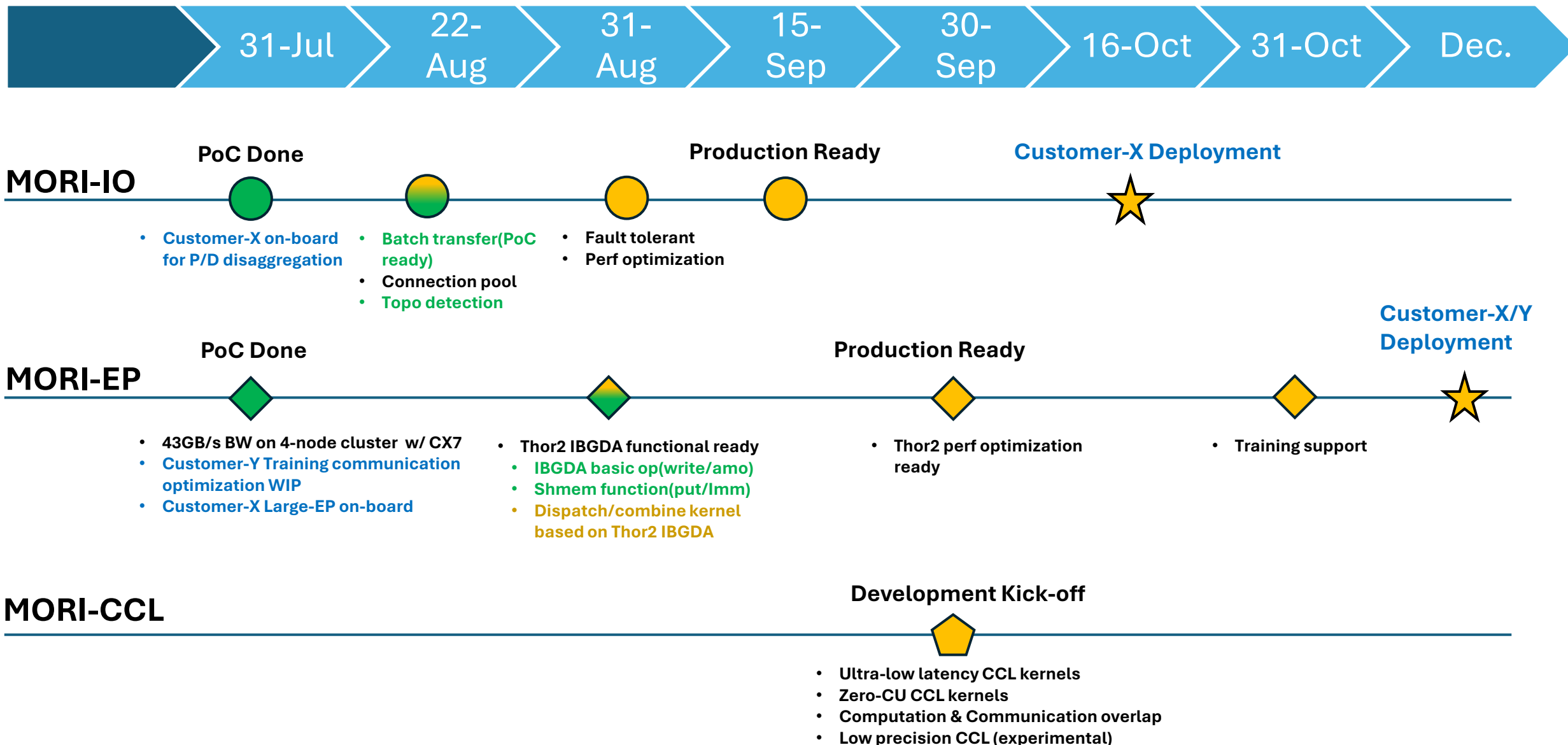
AMD Competitive Landscape



VS.

Nvidia





SGLANG+MORI Integration Evaluating the Tech-Path now!

AMD's LLM Optimization

