## **LLMIR System Architecture**

## **Application Layer (vLLM / SGLang) MLIR Optimization Pipeline Frontend Converters** • PyTorch Models • vLLM Integration • SGLang Integration LLMIR Compiler Core **LLM Dialect Optimization Passes** • PagedKVCacheType KV Cache Fusion $\bullet \ Sharded Tensor Type \\$ • Quantization-Aware • QuantizedTensorType • Parallelization • llm.paged\_attention Memory Optimization **Backend Code Generators** CUDA Backend ROCm Backend CPU Backend **Execution Layer (CUDA / ROCm / CPU)**