



# Buddy Compiler



## An End-to-End AI Compiler from DSL to DSA

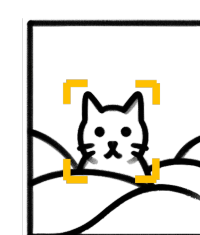
### Buddy Compiler Overview

Homepage: <https://buddy-compiler.github.io/>  
GitHub: <https://github.com/buddy-compiler>

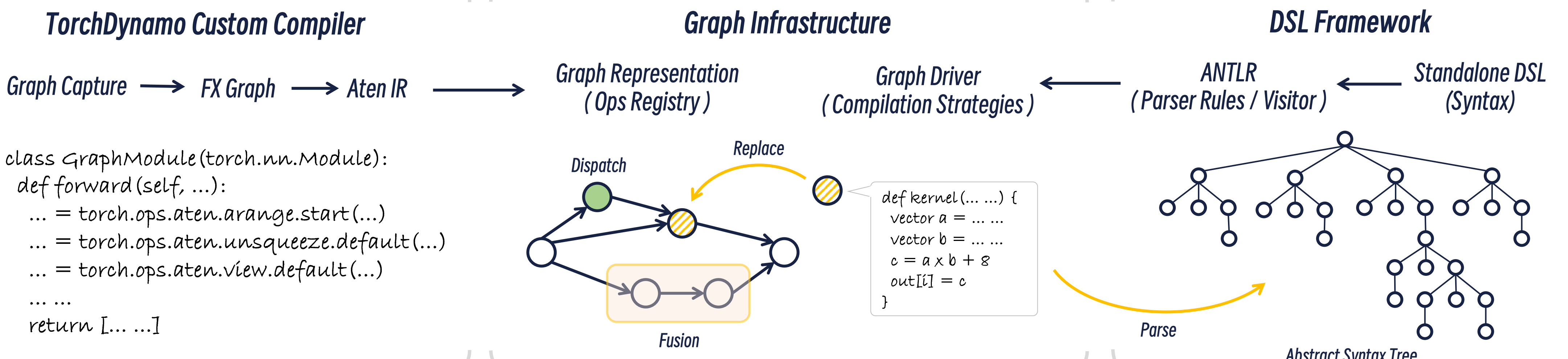
Buddy Compiler is a “buddy system” to bridge the gap between AI models and hardware platforms. It consists of a frontend for AI framework and language integration, a middle-section leveraging MLIR for efficient optimization and tuning, and a backend for advanced code generation and hardware collaboration. By providing an end-to-end pipeline from domain-specific language (DSL) to domain-specific architecture (DSA), Buddy Compiler aims to establish a collaborative AI compiler ecosystem<sup>[1]</sup> to address fragmentation challenges in the current AI software stack.

Buddy Compiler’s frontend seamlessly connects with the DSL framework and PyTorch 2.x compilation framework via a graph infrastructure. The frontend integrates PyTorch AI models and computation kernels written in DSL into a unified graph representation. Moreover, the graph driver performs graph-level transformations and determines the compilation strategy, effectively converting the graph representation into MLIR.

### Buddy Compiler Frontend

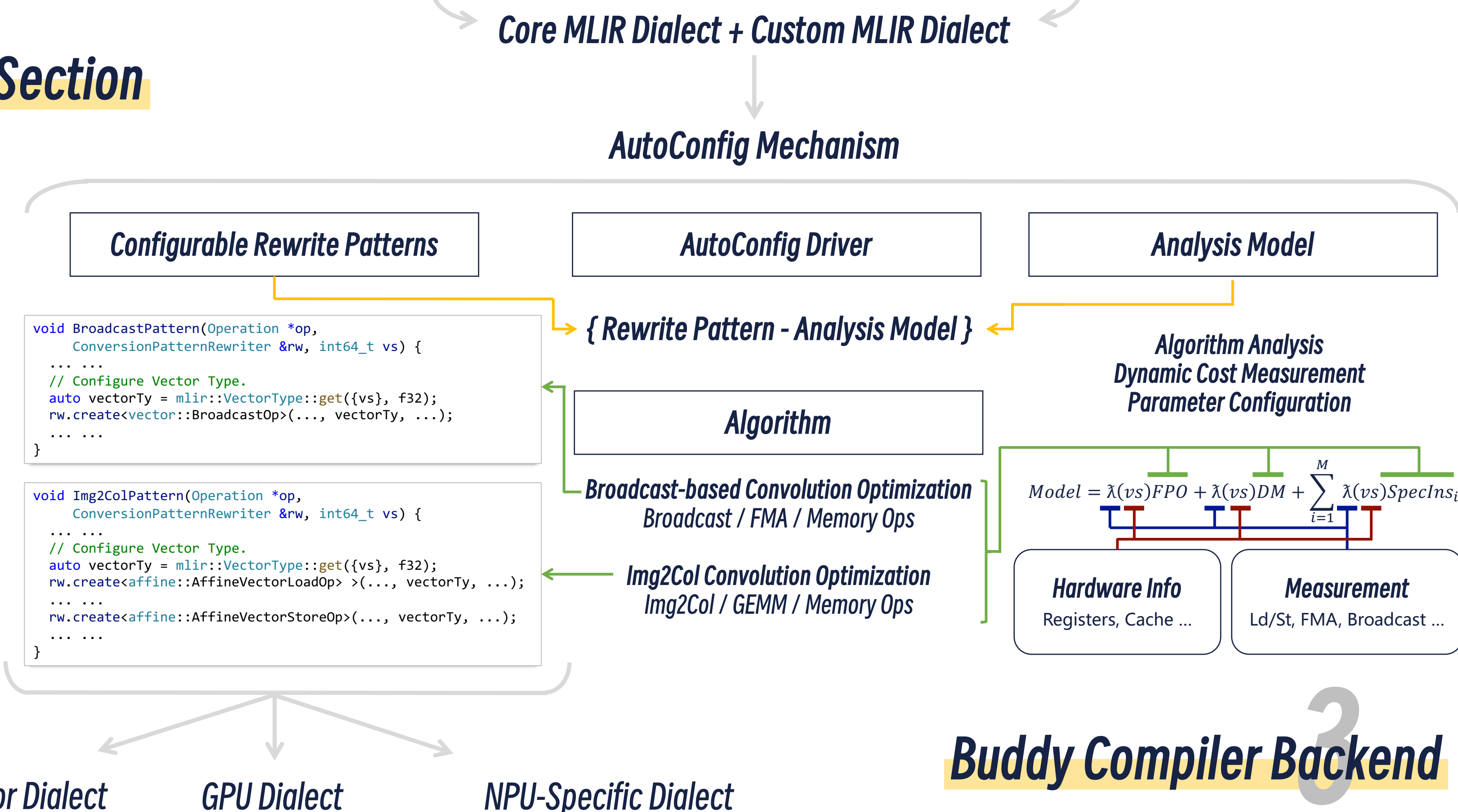


Multimodal Input  
Text / Audio / Image



### Buddy Compiler Middle-Section

Buddy Compiler’s middle-section focuses on performance optimization and domain-specific MLIR extensions (digital audio processing dialect and digital image processing dialect) to support efficient multimodal AI computation. It implements configurable rewrite patterns as compilation optimization strategies and builds the AutoConfig<sup>[2]</sup> mechanism for performance tuning. This approach ensures that optimizations are implemented once, delivering consistent impacts across different platforms.



### Buddy Compiler Backend

Buddy Compiler’s backend is dedicated to hardware-specific code generation and collaborative design, fully leveraging the acceleration capabilities of hardware architectures. When generic abstractions are insufficient to encapsulate the hardware-specific traits, it creates customized representations and compilation pipelines to enable specific optimization strategies.

MLIR HW Dialects:	Extended Vector Dialect	GPU Dialect	NPU-Specific Dialect
System Tools:	LLVM Toolchain + Hardware-Specific Toolchain / Toolkit / Library + Execution Engine		
Target Platforms:	SIMD / Vector Processors (X86 AVX / Arm Neon / RVV)	GPU Platforms (Nvidia GPUs)	NPU Platforms (Gemmini Accelerator)

[1] Hongbin Zhang et al.: Compiler Technologies in Deep Learning Co-Design: A Survey

[2] Hongbin Zhang et al.: AutoConfig: A Configuration-Driven Automated Mechanism for Deep Learning Compilation Optimization

#### Authors & Contributors:

Hongbin Zhang, Xulin Zhou, Liutong Han, Taiqi Zheng, Hongyu Lin, Hanghang Cao, Meng Li, Weijia Li, Jiongjia Lu, Mingjie Xing, Yanjun Wu (ISCAS)

Jiuyang Liu, Zikang Liu, Linquan Wei, Yuliang Li, Zhongyu Qin, Sen Yang, Prathamesh Tagore (PLCT Internship)

ISCAS

中国科学院大学