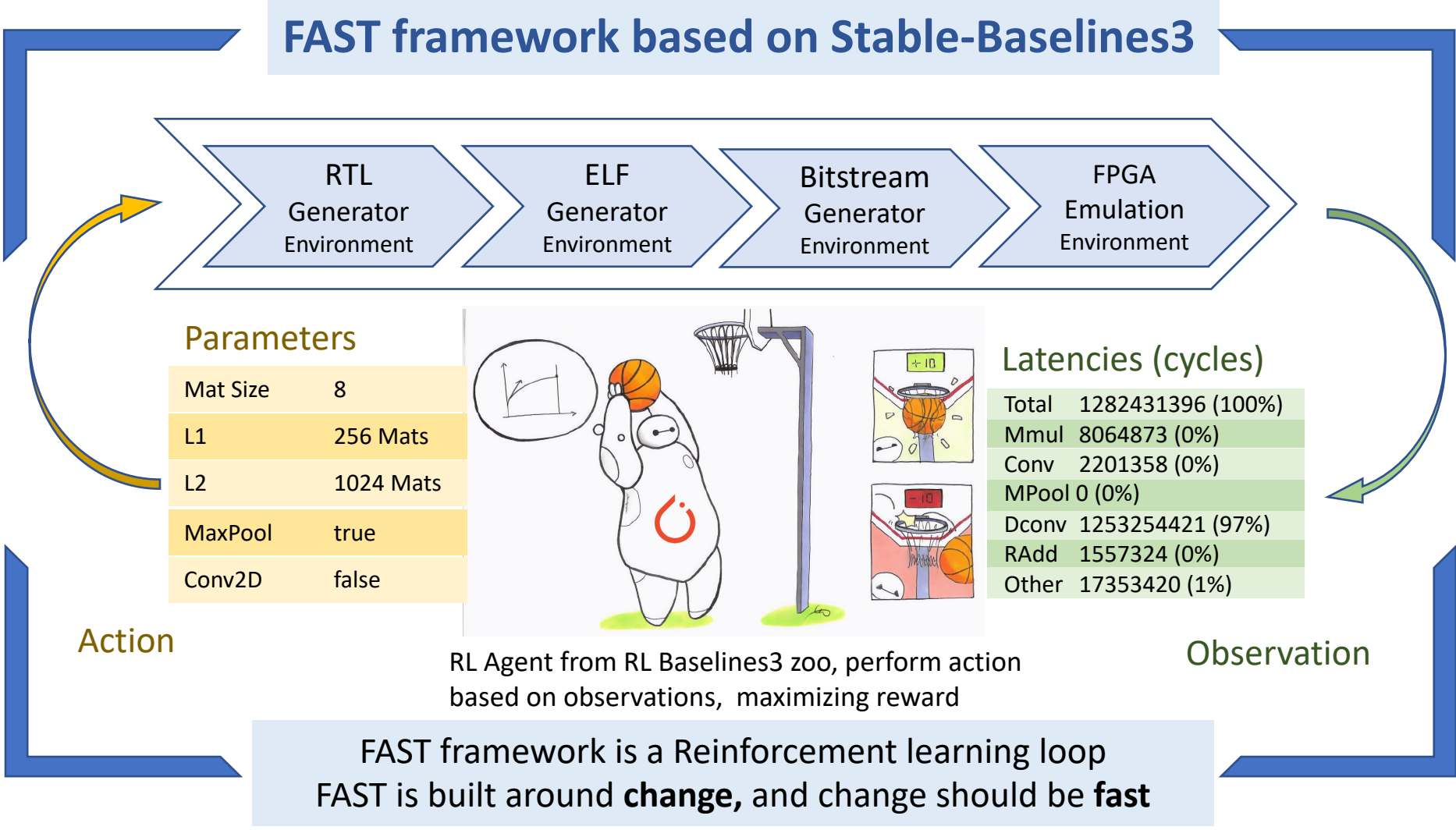


FAST: Accelerating Full-stack SoC Design-space Exploration With FPGA-in-the-Loop



Tayyeb Mahmood, Jaeyong Chung, Incheon National University, Incheon, South Korea

FAST is a work in progress concerning the Architectural and RTL Design-Space Exploration with simulation-driven ML-around approach, connecting the two ends of AIEDA’s bigger picture. Open-source parameterized RTL generators like UC Berkeley’s **Chipyard** and UofU’s **openFPGA** have seen commercial success. To facilitate the optimization of a large design space unleashed by these projects with ML/RL- driven approaches, we present an **FPGA-Assisted Scalable Test** framework. FAST accelerates DSE by emulating the generated RTL and by monitoring the real-time performance and power with FPGA-in-the-Loop, to achieve an accuracy compared to conventional **model-based estimators**, at speeds many orders higher than **cycle-accurate simulators**. The key challenge to be addressed in this research is to reduce the time for RTL-to-Bitstream flow. Utilizing FAST framework, domain-experts can optimize a full-stack of AI-accelerators and reconfigurable arrays at prototyping stage or can search for an optimum design for FPGA deployments, like **Microsoft Brainwave**, and **Amazon EC2-F1** clusters.



FPGA Platform

SingleE V7 Logic Module
Xilinx XC7V2000T-FLG1925
2.5D SSI, 4 SLRs (24 Clock Regions)
RS232 and SDCard Add-on module
8GB DDR3 SODIMM

Labels on the board include: RS232 Console, SDCard Benchmark Workload, JTAG Configuration Debug Reconfig, SysCon Clock and Reset Power Monitor, and XC7V2000T.

FAST Software execution order

On-chip ROM, NorFlash, NANDFlash, On-chip SRAM, DRAM

M-Mode, S-Mode

Rocket-chip oe-linux boot (From left to right)

ZSBL, FSBL, openSBI, u-boot, linux

FAST boot: FSBL, Workload

SDCard partition map: meta-sifive GPT partitions (fsbl, u-boot, ext4, user), FAST MBR partition (ARM/MLO style) (fat32)

Stable-Baselines3 framework

A set of reliable implementations for Reinforcement Learning, like DQN, PPO, etc.
A rework of OpenAI baselines with a number of enhancements, based on Gym.

OpenAI Gym
Provides a standard API to communicate between RL agents and environments, and utilities to create, analyze, evaluate and vectorize the environments, action spaces and observation spaces.

RL Baselines3 Zoo
Provides framework for training and evaluating RL agents, and tuning hyperparameters
Provides a collection of tuned hyperparameters for common environments and RL algorithms, and agents.

Gym Space
A data structure to support multi-dimensional observation and action spaces, like discrete, continuous, maps and graphs.
Support random sampling of actions and validations of observations.

Gym Environment
A wrapper around a Markov decision process to learn, e.g. a retro game, a traffic controller or parameter tuner as FAST.
Supports two simple functions, reset and step.

Gym Vectorization
Provides framework to run multiple instances of Environments, in parallel (multiprocessing).
Input: a batch of actions, Output: a batch of observations
Scalability in clusters and cloud.

FAST framework

Goal
To automate, decouple and accelerate the DSE flow from Architectural parameters to RTL to Bitstream to Execution, as shown in the figure.

Challenge
RTL to bitstream flow is optimized for a time-consuming single pass, followed by iterative optimizations. Incremental design flow only allows minor touching.

Motivation
DSE flow can take advantage of hierarchical design, with synthesis and P&R of only those modules which are modified by parameters.

RTL Test Harness- Chipyard

SDCard, UART, BSCAN (Debug), Rocket-chip, L2 Cache, Gemmini, BSCAN (Reconfig), Clock & Reset, ChipTop, Shell

Shell: Sifive’s fpga-shells framework
BSCAN: Xilinx Primitive instance
DDRCon: Xilinx DDR3 MIG
SDCard: Sifive’s sifive-blocks SPI
UART: Sifive’s sifive-blocks UART
Rocket-chip: UCBAR’s RISCv core
L2Cache: Sifive’s inclusive L2 cache
Gemmini: UCBARs systolic array for MatMul, accelerator units for conv, relu and maxpool.
Clock & Reset: 200MHz master clock and a global reset from SysCon

We divide time-consuming RTL-to-Bitstream pass into a sequence of decoupled subprocesses, in order to optimize each subprocess for DSE, individually.

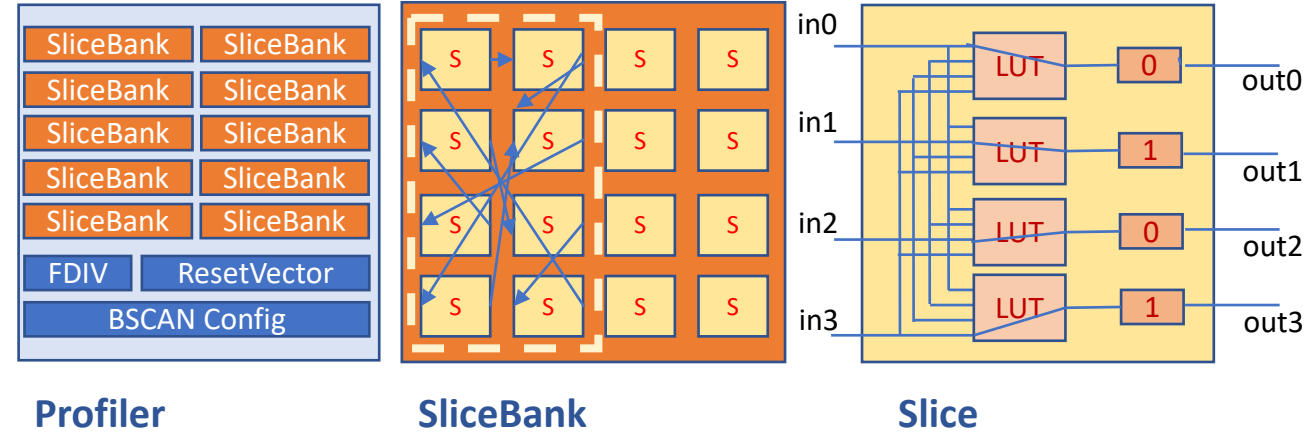
FAST Track

RTL generator
ChipyardFIRRTLEnv (Scala project is assembled into standalone jar to reduce FIRRTL generation time and for scalability.)
ChipyardVerilogEnv (Verilog generation time can be reduced by differential conversion.)

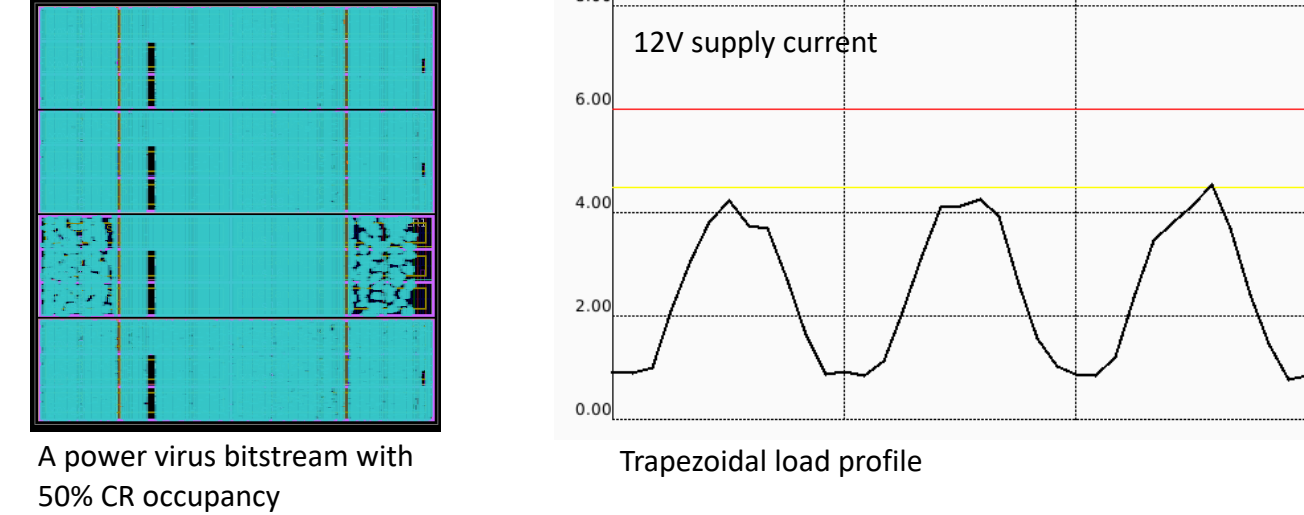
ELF generator
WorkloadGenEnv (c/c++ Compilation is fast)
WorkloadDeployEnv (ELFs must be prestored on SDCard, currently it is manual process.)

Bitstream generator
FPGASynthesisEnv (Full synthesis consumes more than 30 minutes, partial synthesis consumes 10 minutes.)
FPGAPlaceEnv (Full Placement consumes 20 minutes, partial placement takes 10 minutes.)
FPGARouteEnv (Full pouting consumes more than 40 minutes, partial routing consumes 10 minutes.)
BitstreamGenEnv (Full bitstream generation consumes more than 5 minutes, partial bitstream consumes < 1 minute.)

FPGA emulation
BitstreamConfigEnv (Configuration is fast.)
FPGAEmulationEnv (Less than 7 minutes for Resnet50, however most of time is consumed in loading image from SDCard.)
ObservationEnv (Scanning data from logs is fast.)



Profiler has slicebanks, one for each clock region.
FDIV is a counter for divided clocks, one of which is MUXed out. ResetVector is a shift register, clocked with FDIV output. Each bit of ResetVector drives a Slicebank reset.
BSCAN config is a capture register on FPGA JTAG tunnel. Config register drives FDIV MUX and ResetVector pattern. SliceBank has a number of slice groups.
A slice group is a group of slices, the even outputs of whom are randomly connected to the odd inputs, and vice versa. Slice models an FPGA slice, the registers of which are reseted to ‘0101....’ pattern, each LUT passing one of input to a register.
The even-odd connections ensure that the internal node have maximum activity at each clock edge, and randomization ensures the use of routing resources, to maximize power consumption.



Conclusion

We note that partial reconfiguration significantly reduces Bitstream generation time, however total test time still exceeds 30 minutes. Although it is an order better than cycle-accurate simulators for designs of this complexity, we can further reduce this time with techniques such as differential generation and synthesis, using a netlist cache and using RapidWright framework for data driven placement and routing. FAST environments can be used for DSE and policy optimization where a policy may be a sequence of parameters with minimum changes at each step, like Gray coding.

