



SHAPING THE NEXT GENERATION OF ELECTRONICS

JUNE 23-27, 2024

MOSCONE WEST CENTER
SAN FRANCISCO, CA, USA





SHAPING THE NEXT GENERATION OF ELECTRONICS

JUNE 23-27, 2024

MOSCONE WEST CENTER
SAN FRANCISCO, CA, USA



Knowing The Spec to Explore The Design via Transformed Bayesian Optimization

Donger Luo^{*1}, Qi Sun^{*2}, Xinheng Li¹, Chen Bai³, Bei Yu³, Hao Geng^{1,4}

¹ShanghaiTech University ²Zhejiang University

³Chinese University of Hong Kong

⁴Shanghai Engineering Research Center of
Energy Efficient and Custom AI IC

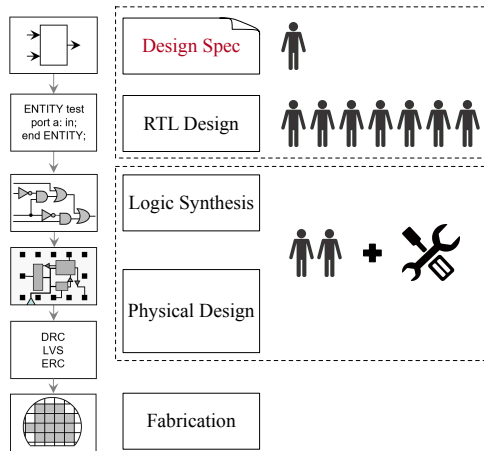


Outline

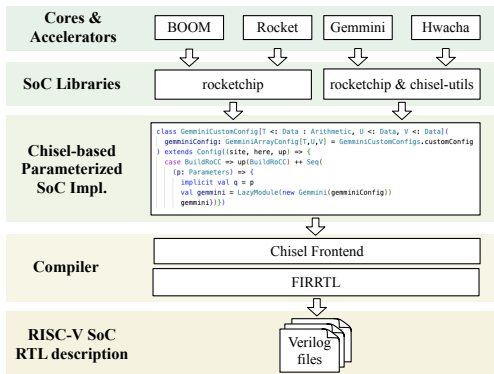
- 1 Introduction
- 2 Algorithms
- 3 Experiment
- 4 Conclusion

Digital IC Design Flow

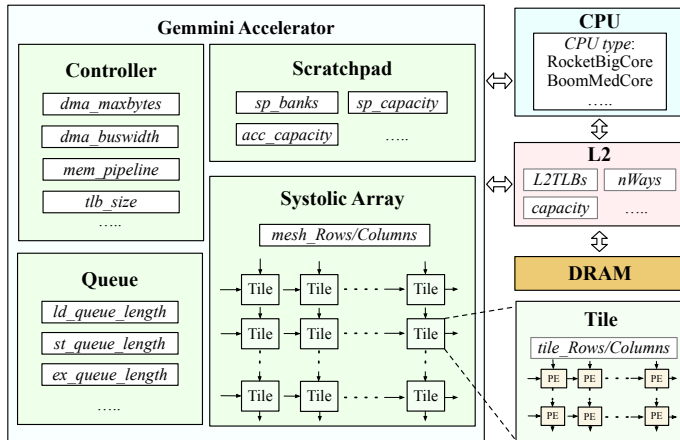
- Digital IC design back end is automated using powerful EDA tools.
- While the front end of digital IC design still requires lots of manpower.



Agile Design Based on Chisel



- Agile development is gradually being adopted to reduce chip design costs and accelerate design cycles.
- Chisel designs, like Rocket Core, Berkeley Out-of-Order Machine, and Gemmini, are configurable and parameterizable RISC-V processors.



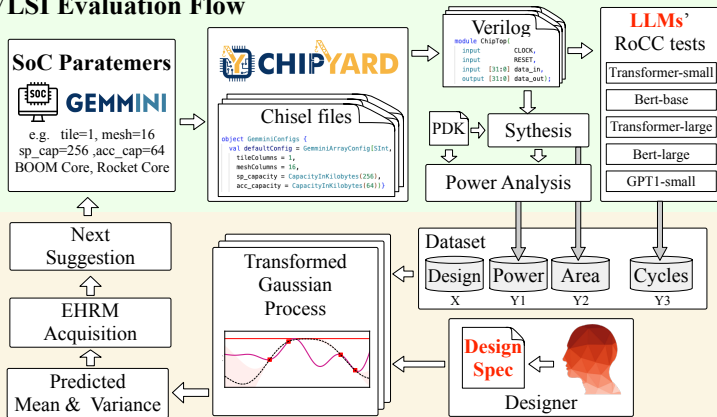
Gemmini SoC Parameters Examples

Parameters (19 in the paper) include:

- selection of CPU
- cache sizes
- accelerator configurations

Parameters	Stage	Candidates
<i>cpu_type</i>	CPU core	RocketBig/BoomMed/BoomLarge
<i>L2TLBs</i>	L2 Cache	512, 1024
<i>nWays</i>		4, 8, 16
<i>capacityKB</i>		512, 1024
<i>tile_Rows/Columns</i>	Accelerator	1, 4, 8
<i>mesh_Rows/Columns</i>		8, 16, 32, 64
<i>sp_capacity</i>		256, 512, 1024, 2048, 4096
<i>sp_banks</i>		4, 16
<i>dma_buswidth</i>		128, 256

VLSI Evaluation Flow



Transformed Bayesian Optimization

Vanilla Gaussian Process

- A GP model is defined as a collection of random variables, any finite number of which have a joint Gaussian distribution.
- A GP is completely specified by its mean function $m(\mathbf{x})$ and covariance function $k(\mathbf{x}, \mathbf{x}')$:

$$f(\mathbf{x}) \sim \mathcal{GP}(m(\mathbf{x}), k(\mathbf{x}, \mathbf{x}')), \quad (1)$$

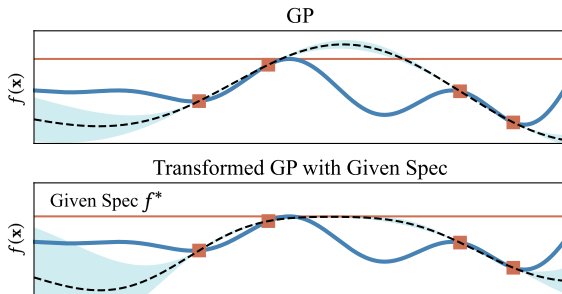
where $\mathbf{x} \in \mathbb{R}^d$ represents the input variable vector.

Transformed Gaussian Process

The Gaussian process is transformed as:

$$f(\mathbf{x}) = f^* - \frac{1}{2}g^2(\mathbf{x}) \quad g(\mathbf{x}) \sim GP(m_0, K), \quad (2)$$

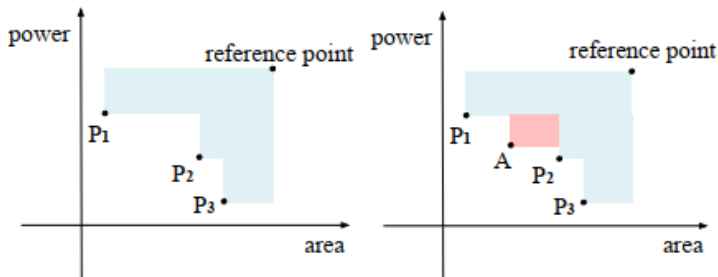
so $f(x)$ will not be beyond the given spec's value f^* .



Acquisition Function

The expected hypervolume improvement EHVI is defined as the expectation of hypervolume's improvement with respect to the posterior predictive distribution of the GP:

$$\text{EHVI}(\mathbf{y}) = \mathbb{E}_{p(\mathbf{y}|D)}[I(\mathbf{y})], \quad (3)$$



EHRM aims to find a configuration of SoC architecture parameters with an expectation closet to the given spec, in other words, with the smallest expected hyper-regret:

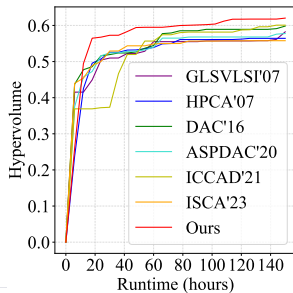
$$\mathbf{x}_{t+1} = \arg \min_{\mathbf{x} \in X} \text{EHRM}(\mathbf{x}) = \arg \min_{\mathbf{x} \in X} \mathbb{E}[Hr(\mathbf{x})]. \quad (4)$$

Experimental Setting

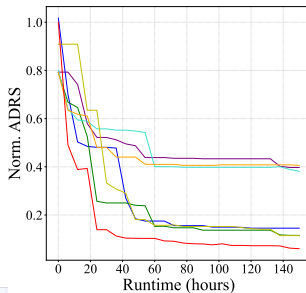
- Gemmini-based RISC-V SoC
- 5 Popular LLMs as performance evaluation
- Cadence Joules and Genus for power and area evaluation
- ASAP7 7nm PDK

Experiment

Metric \ Methods	GLSVLSI'07	HPCA'07	DAC'16	ASPDAC'20	ICCAD'21	ISCA'23	Ours
HV _{0,1}	0.6320	0.6491	0.6789	0.6610	0.6716	0.6398	0.7063
HV _{0,2}	0.7129	0.7255	0.7231	0.7144	0.7251	0.6929	0.7472
HV	0.5577	0.5636	0.5891	0.5758	0.5975	0.5609	0.6208
Average	0.6342	0.6460	0.6637	0.6504	0.6647	0.6312	0.6914
Ratio(%)	91.72	93.43	95.99	94.07	96.13	91.29	100



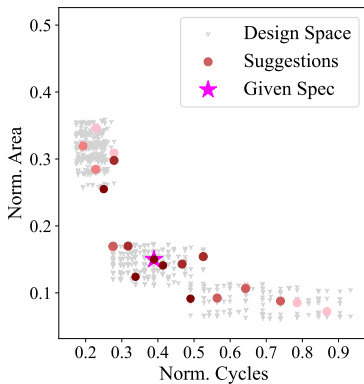
(a)



(b)

Experiment

- The suggestion keeps getting close to the given spec QoR values and finally reaches the targets.



Conclusion

- An architecture design space exploration method based on the transformed Bayesian optimization approach.
- The constructed model utilizes the given spec QoR metric values as additional information to learn.
- A tailored acquisition function is developed for optimization in multiple metrics (e.g., cycles, power, and area).



SHAPING THE NEXT GENERATION OF ELECTRONICS

JUNE 23-27, 2024

MOSCONE WEST CENTER
SAN FRANCISCO, CA, USA



Thanks!

