

# Dynamic Adaptation of Runtime Systems Based on Energy Consumption

IFL'24 – 27 August 2024

**Jordy Aldering**  
Bernard van Gastel  
Sven-Bodo Scholz

# Energy crisis

**Greenhouse gas emissions of ICT projected to be 14% of total in 2040** [1]

**We should take steps to decrease energy consumption of software!**

## **However**

- Measuring energy consumption is hard
  - Cannot isolate a single process
- Decreasing energy consumption is even harder!
  - No control over environment and background

[1] Lotfi Belkhir, Ahmed Elmeligi, *Assessing ICT global emissions footprint: Trends to 2040 & recommendations*

# Adapt to runtime changes

## Cannot isolate energy consumption of a single process

- Consider system's energy consumption as a whole

## What users can measure

- CPU total energy consumption (RAPL)

## What we can control

- Thread-count
- Perfect for HPC
- Interested in the variation between measurements

Why is a static choice not good enough?

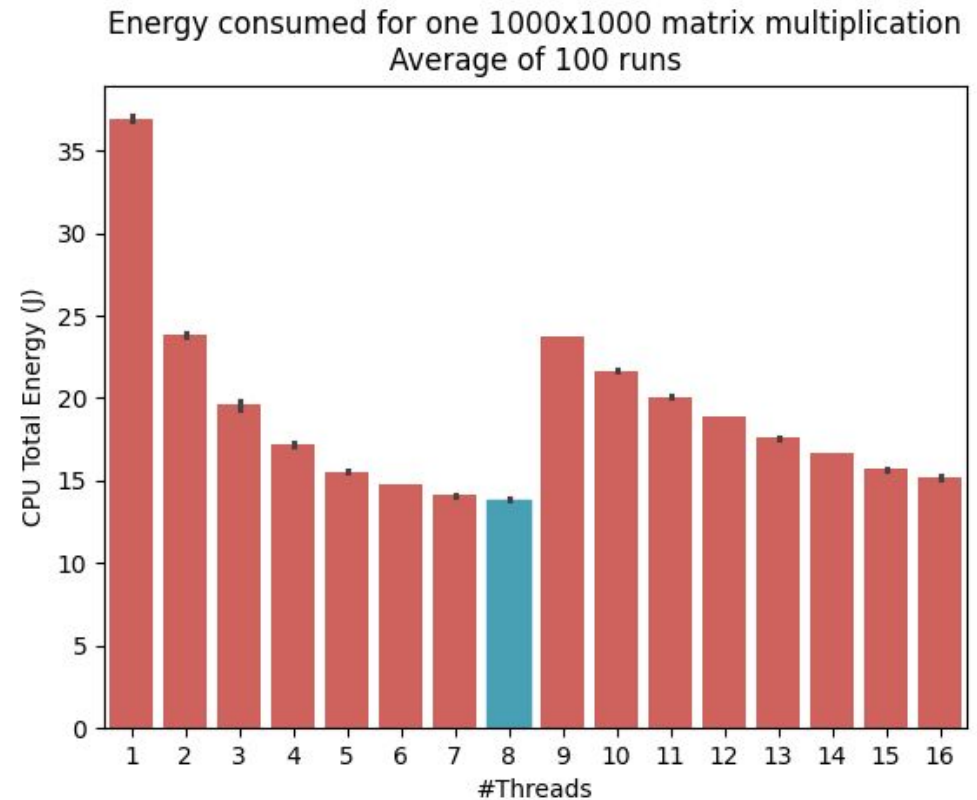
# Energy per thread-count

## Intel Xeon E-2378

- 8 homogeneous cores
- 16 threads

## What about

- Other running programs
- Different hardware
- Performance/efficiency cores
- Power profile
- Architecture
- Throttling
- Pinning
- etc.



Why is a static choice not good enough?

# Energy per thread-count

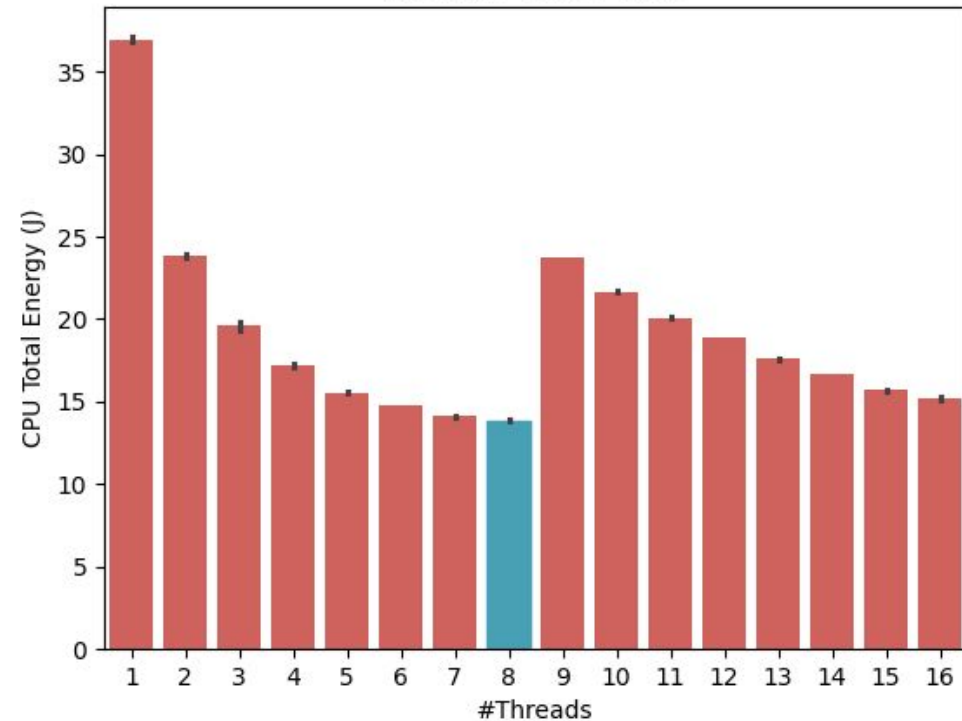
## Intel Xeon E-2378

- 8 homogeneous cores
- 16 threads

## What about

- Other running programs
- Different hardware
- Performance/efficiency cores
- Power profile
- Architecture
- Throttling
- Pinning
- etc.

Energy consumed for one 1000x1000 matrix multiplication  
Average of 100 runs

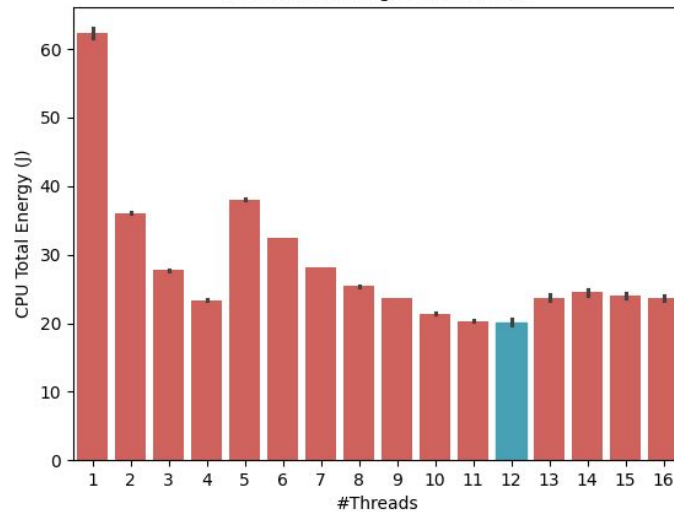


Why is a static choice not good enough?

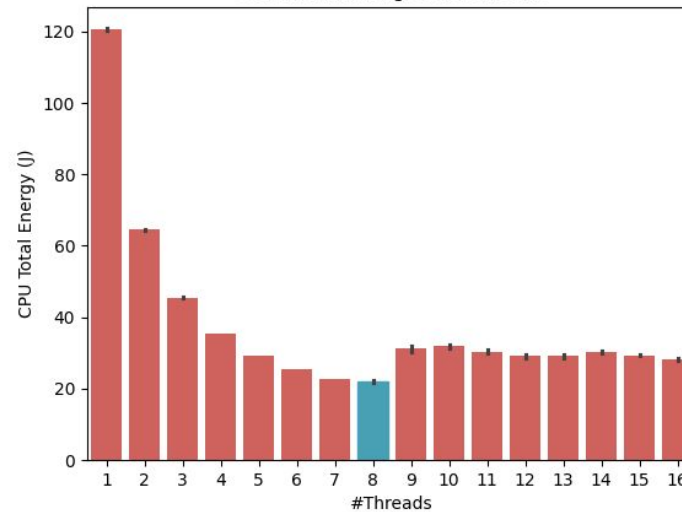
# Energy per thread-count

**How does optimum change if we add a background load?**

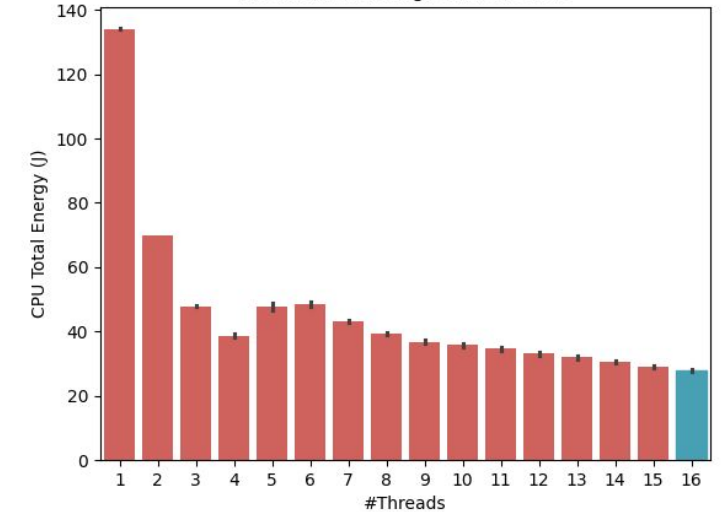
Energy consumed per 1000x1000 matrix multiply (100 run average)  
4 artificial background threads



Energy consumed per 1000x1000 matrix multiply (100 run average)  
8 artificial background threads



Energy consumed per 1000x1000 matrix multiply (100 run average)  
12 artificial background threads



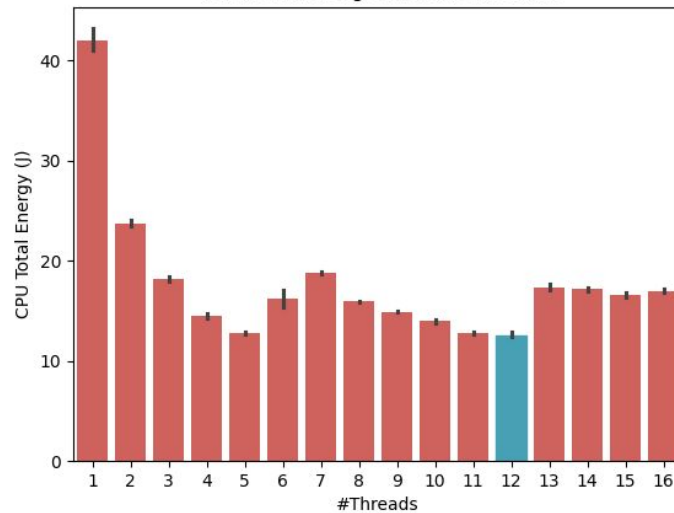
Why is a static choice not good enough?

# Energy per thread-count

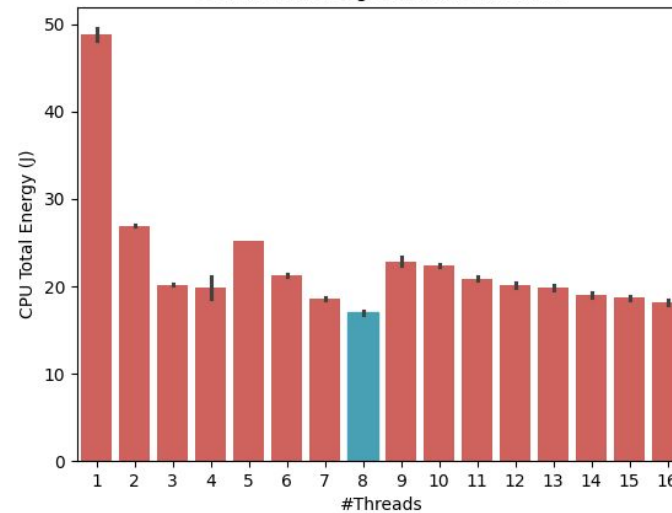
## AMD Ryzen 7 5800H

- 8 homogeneous cores
- 16 threads

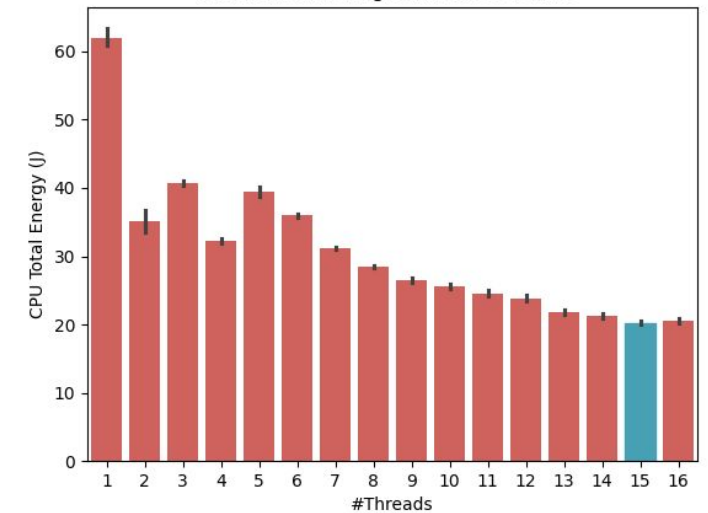
Energy consumed per 1000x1000 matrix multiply (100 run average)  
4 artificial background threads AMD



Energy consumed per 1000x1000 matrix multiply (100 run average)  
8 artificial background threads AMD



Energy consumed per 1000x1000 matrix multiply (100 run average)  
12 artificial background threads AMD



# Dynamic Adaptation

**Shows why we need dynamic adaptation**

**Note that**

- System changes relatively slow compared to measuring time
  - Other relatively long-running programs
  - Otherwise we consider it as noise
- Require a repeated parallel application
  - Matmul
  - Stencils
  - N-body
  - etc.



# How?

## **Program provides energy statistics of the parallel iteration**

- With-loop in the case of SaC
- Energy measured with RAPL

## **Framework updates optimal thread-count every 20 iterations**

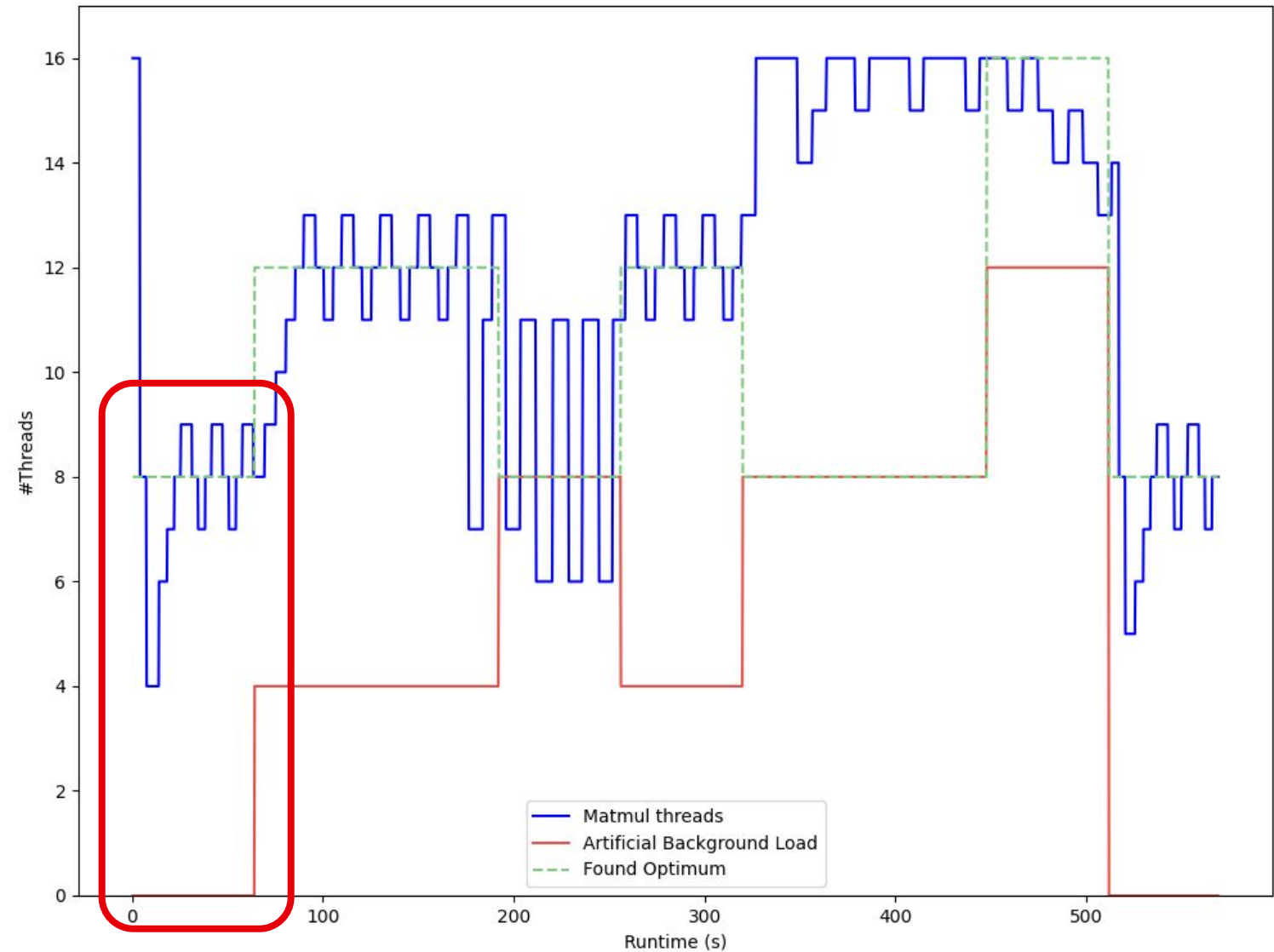
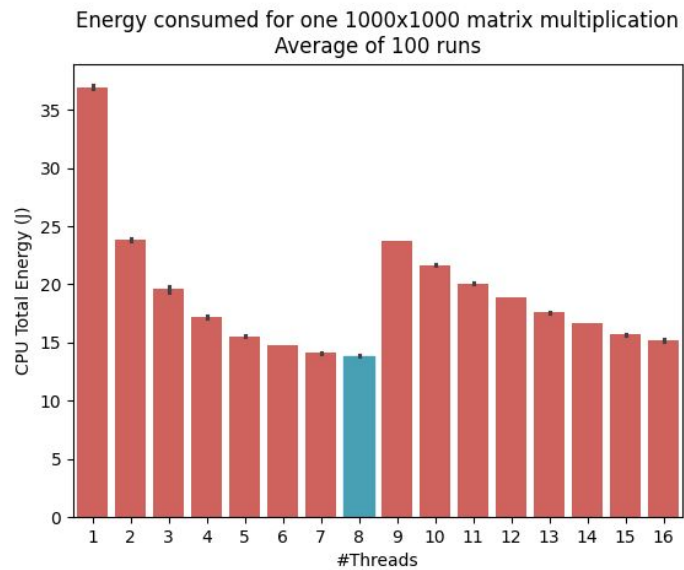
- Get lowest energy consumption those measurements
  - Frequency distribution
- Compare against
  - Previous energy consumption
  - Currently known best energy consumption
- Accordingly adjust step direction and step size

# Results

Figure 10 is a line graph showing the number of threads over runtime for different background loads. The x-axis is 'Runtime (s)' from 0 to 550. The y-axis is '#Threads' from 0 to 16. Three series are plotted: 'Matmul threads' (solid blue line), 'Artificial Background Load' (solid red line), and 'Found Optimum' (dashed green line). The blue line fluctuates between 4 and 16 threads. The red line shows background load steps at 4, 8, and 12 threads. The green dashed line indicates the found optimum, which is 8 threads for low background load, 12 threads for medium background load, and 16 threads for high background load.

Changing background load

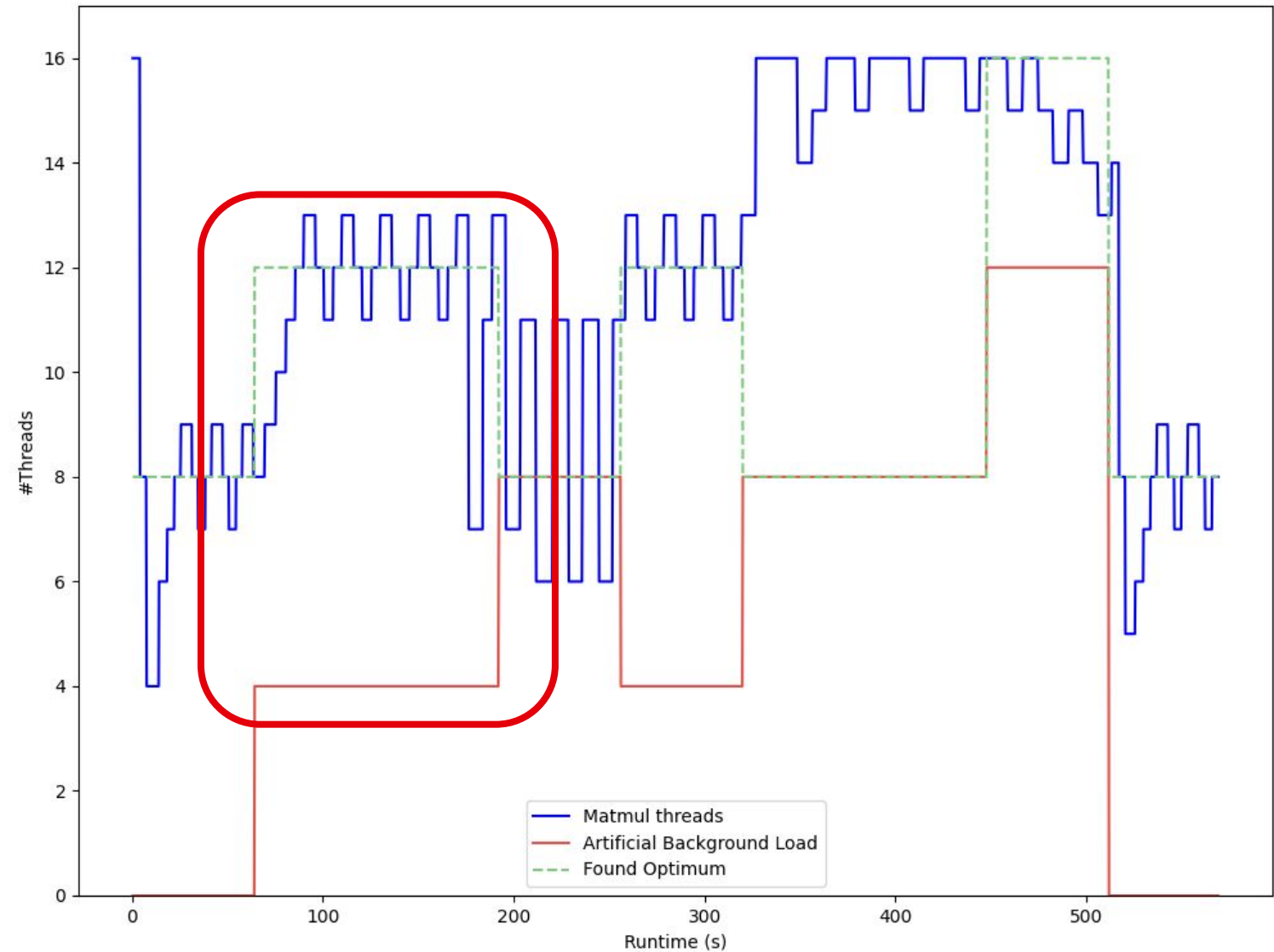
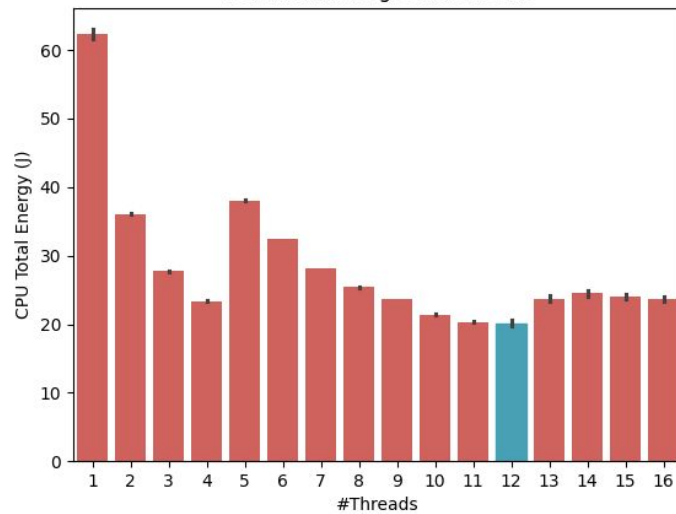
# Results



Changing background load

# Results

Energy consumed per 1000x1000 matrix multiply (100 run average)  
4 artificial background threads

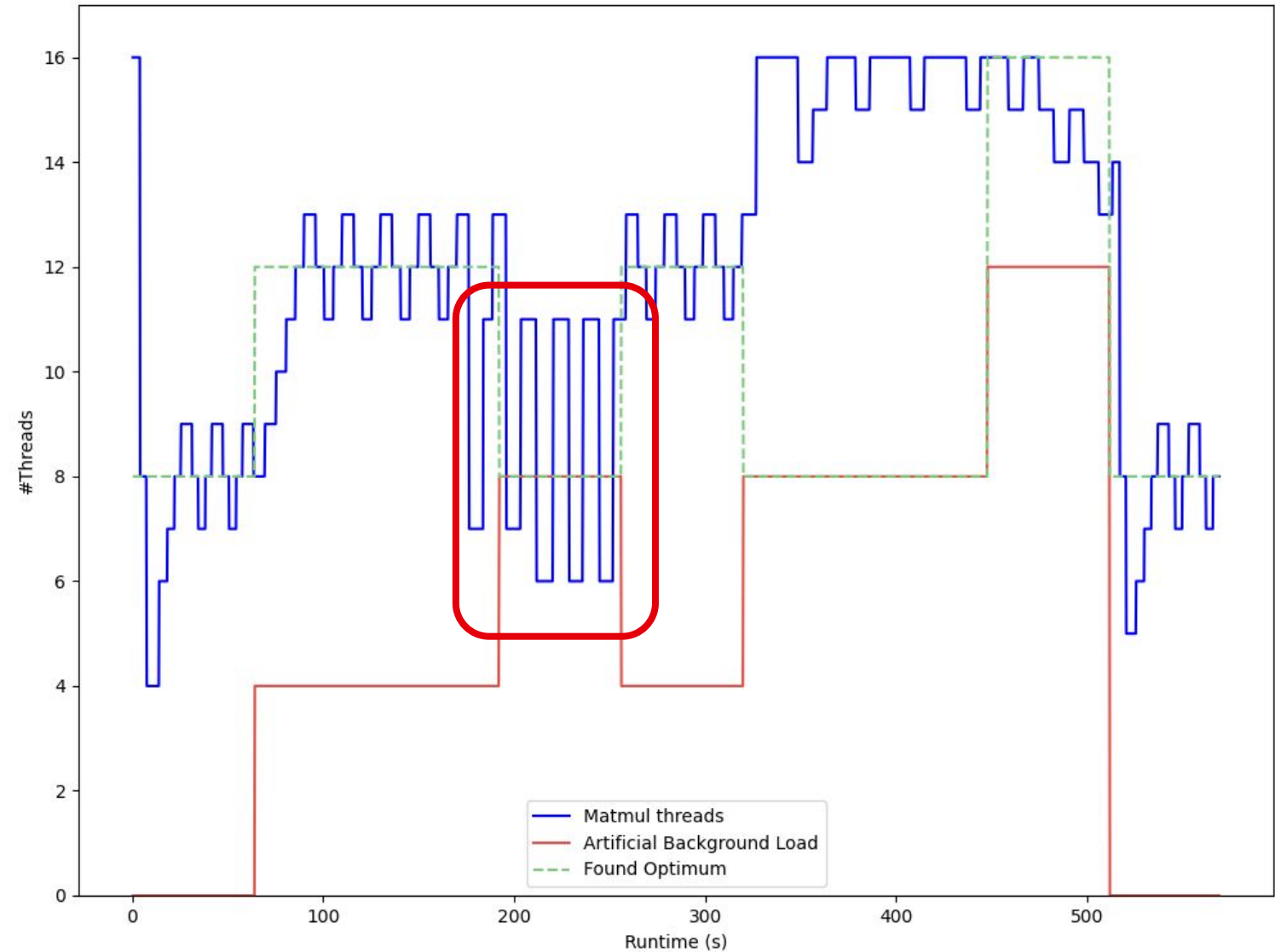
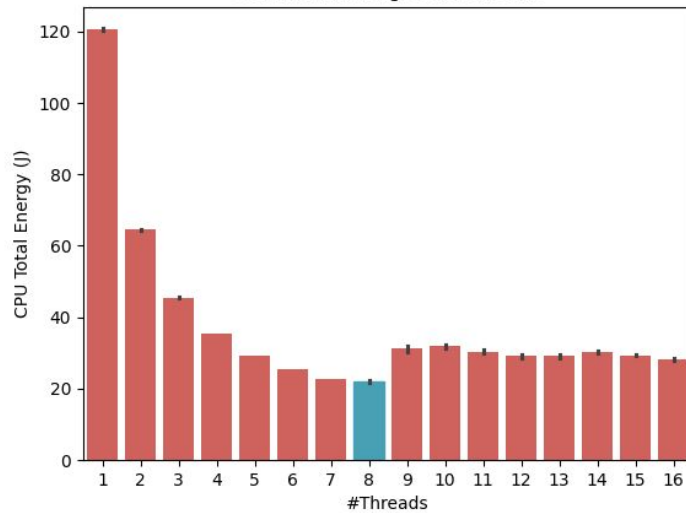


Changing background load

# Results

**We keep jumping too far past the optimum**

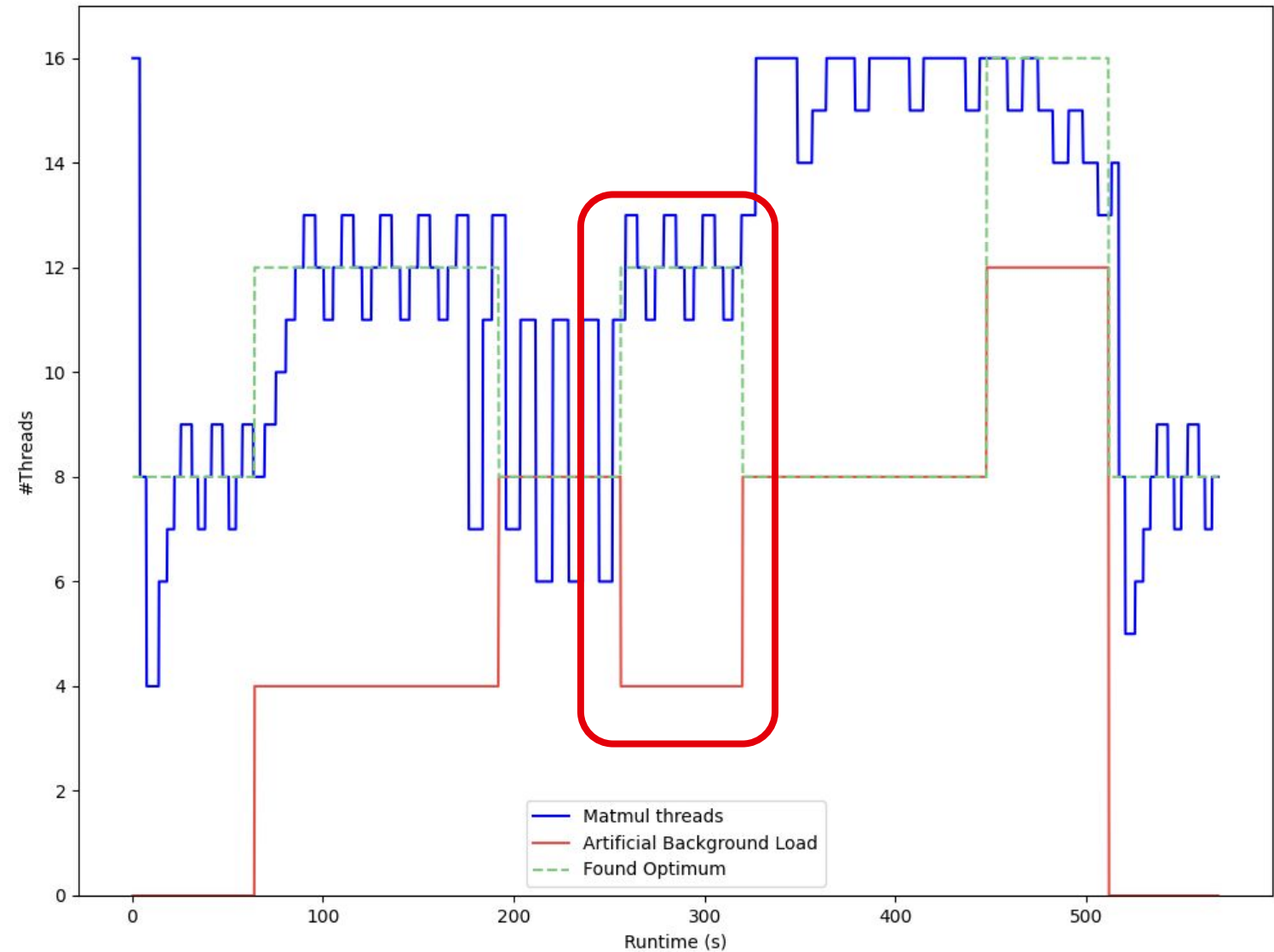
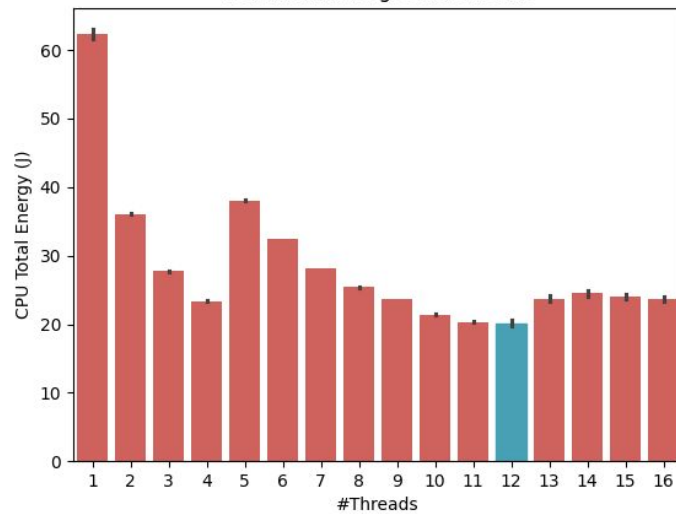
Energy consumed per 1000x1000 matrix multiply (100 run average)  
8 artificial background threads



Changing background load

# Results

Energy consumed per 1000x1000 matrix multiply (100 run average)  
4 artificial background threads

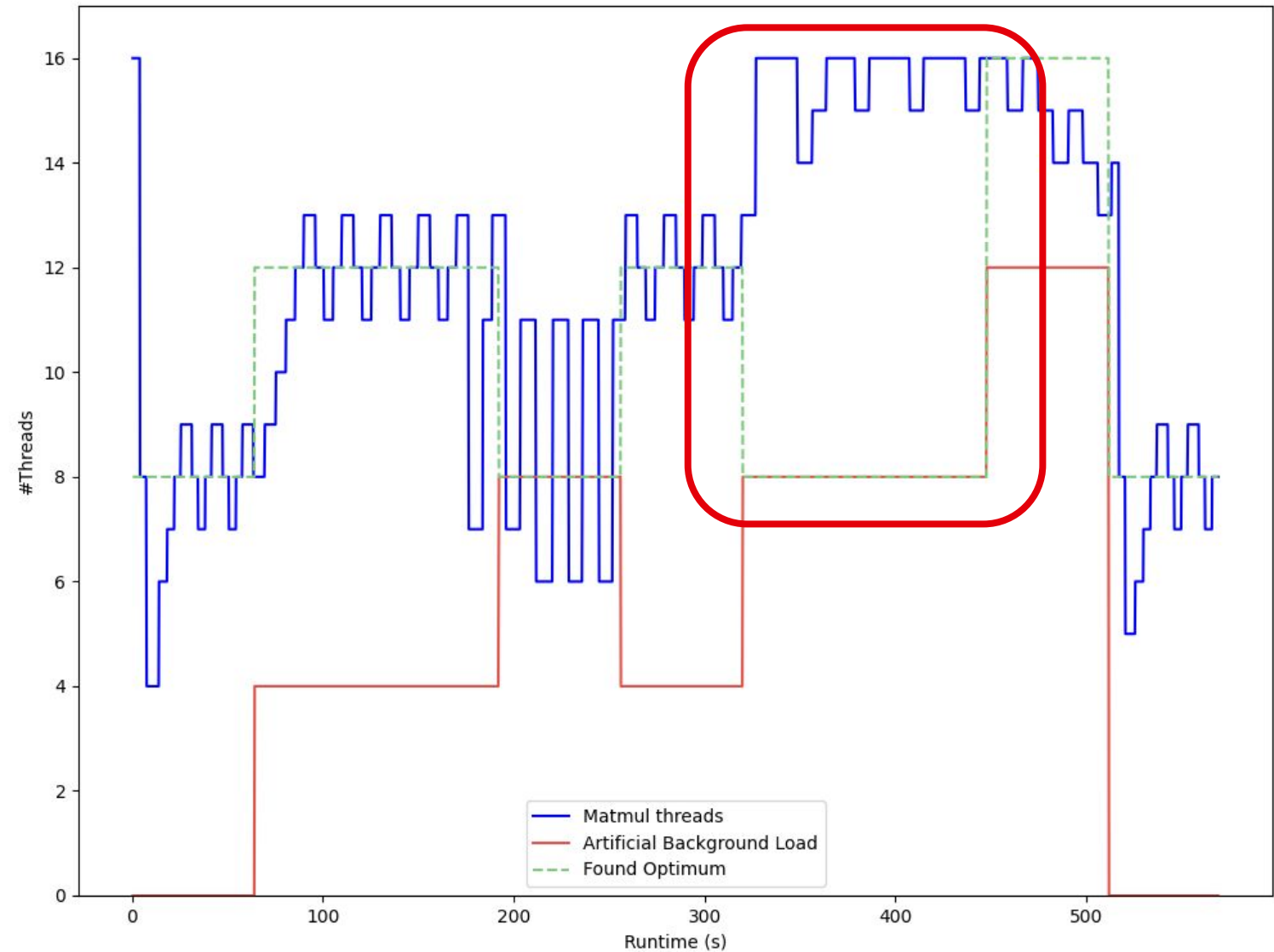
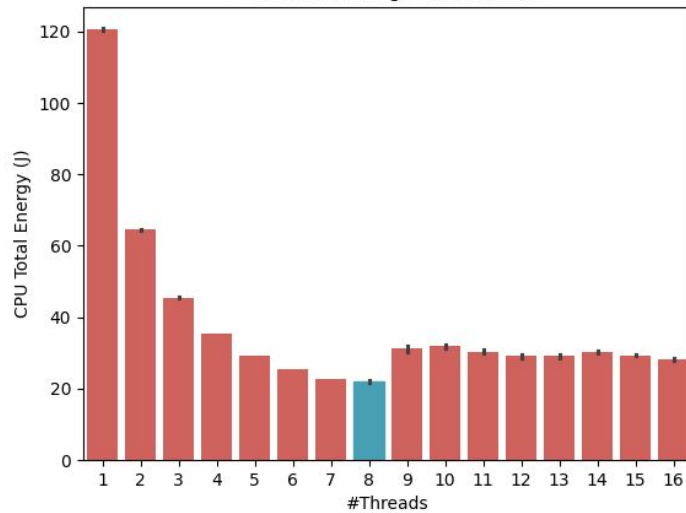


Changing background load

# Results

## Local optimum

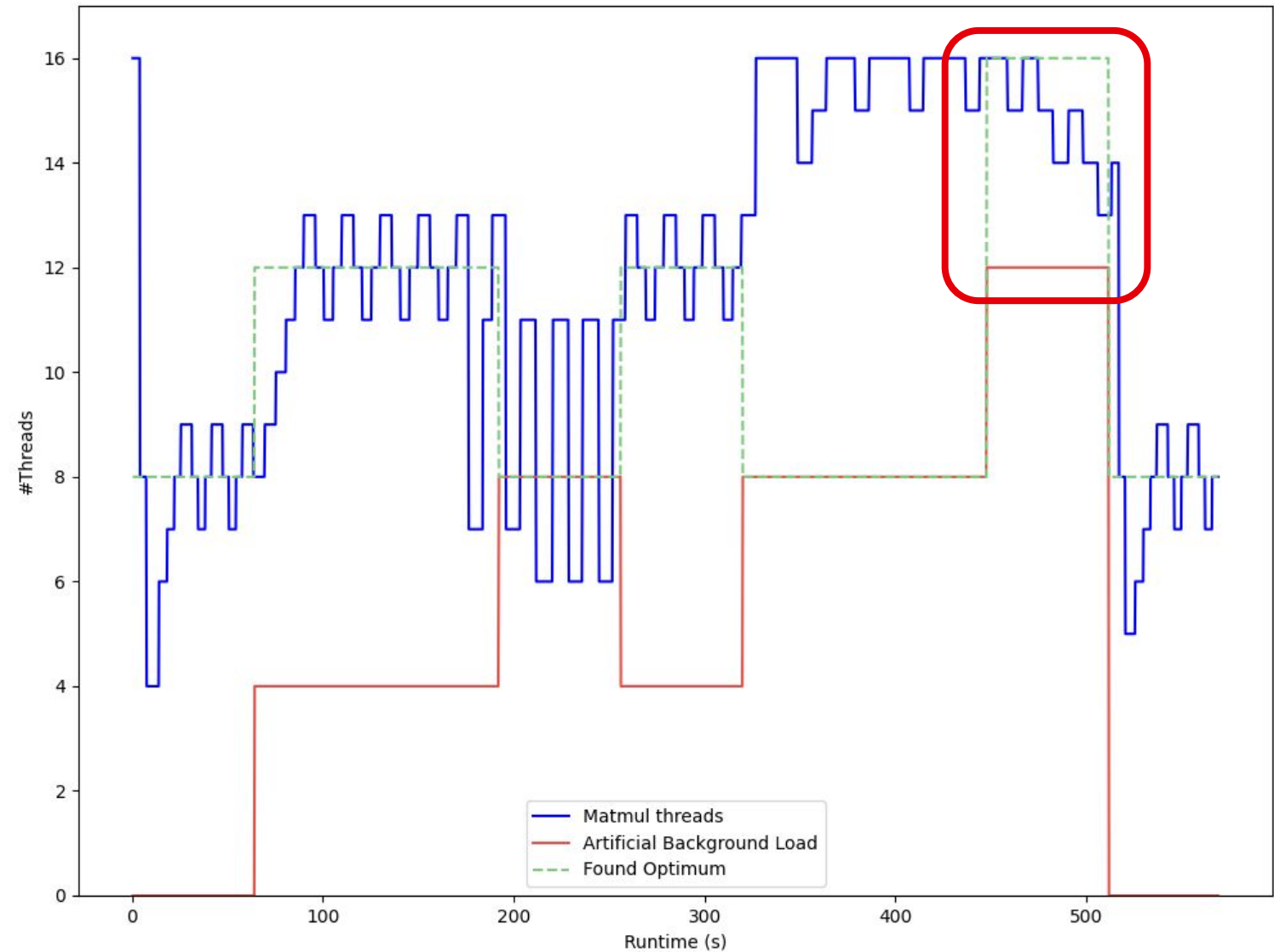
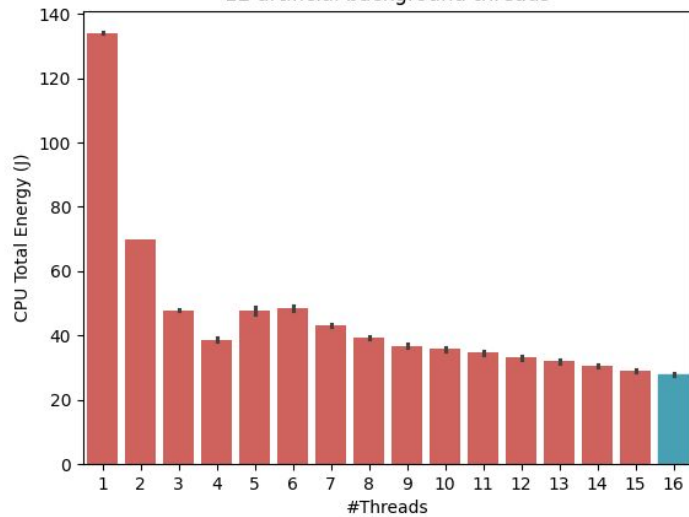
Energy consumed per 1000x1000 matrix multiply (100 run average)  
8 artificial background threads



Changing background load

# Results

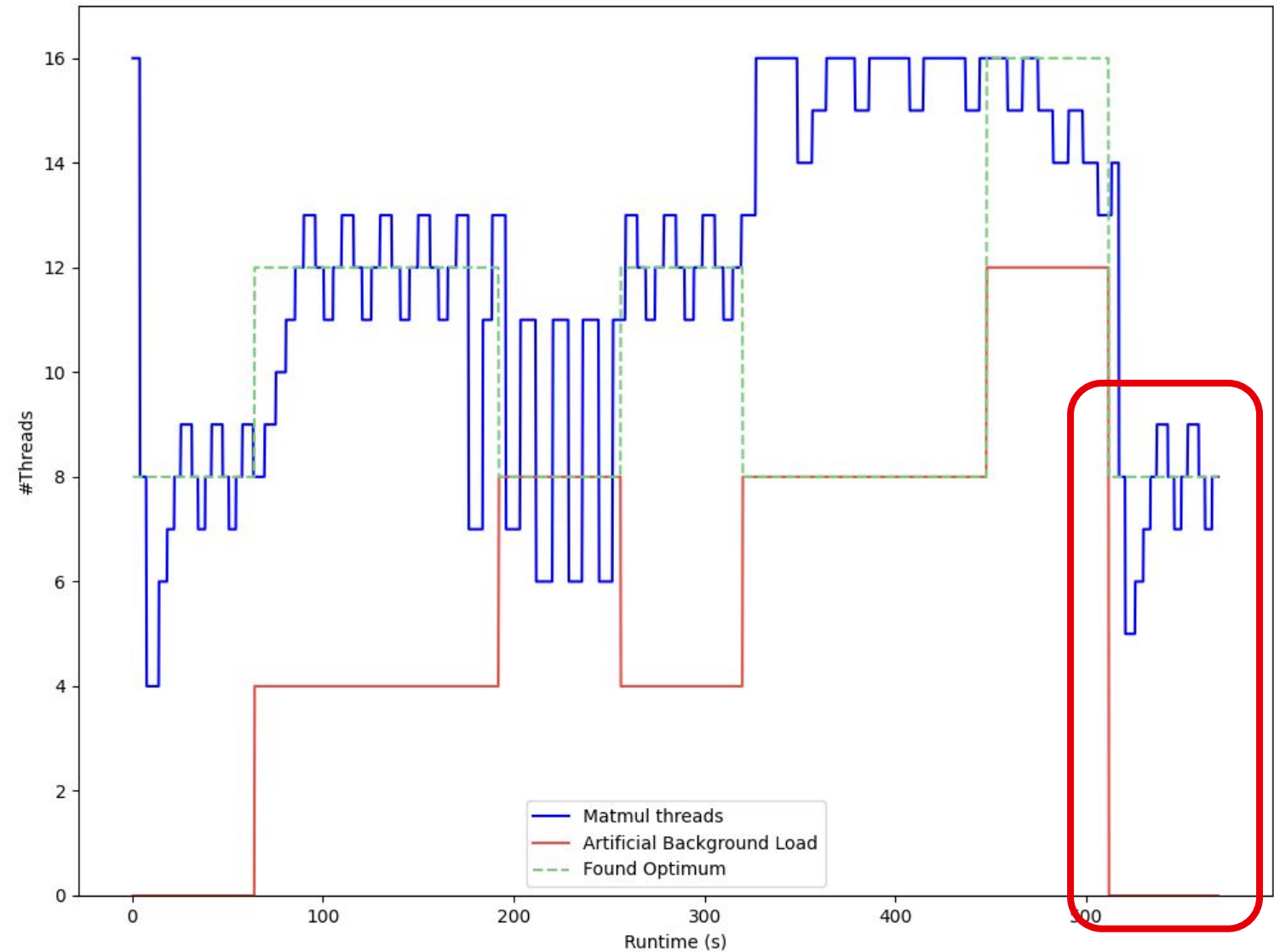
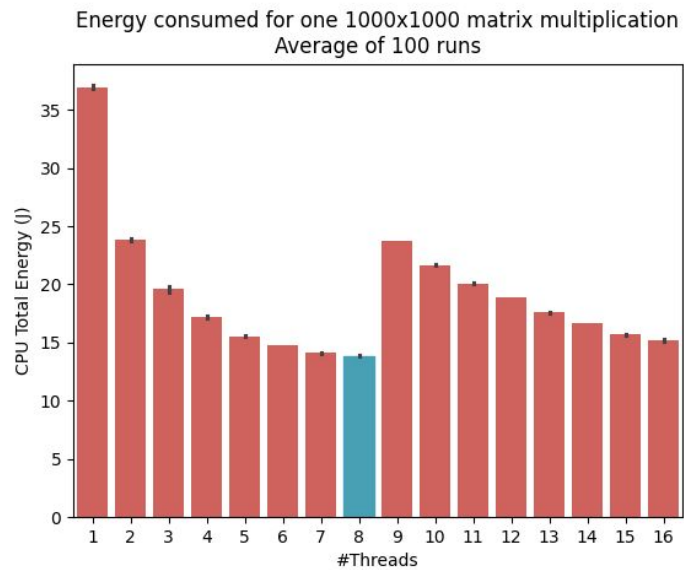
Energy consumed per 1000x1000 matrix multiply (100 run average)  
12 artificial background threads





Changing background load

# Results

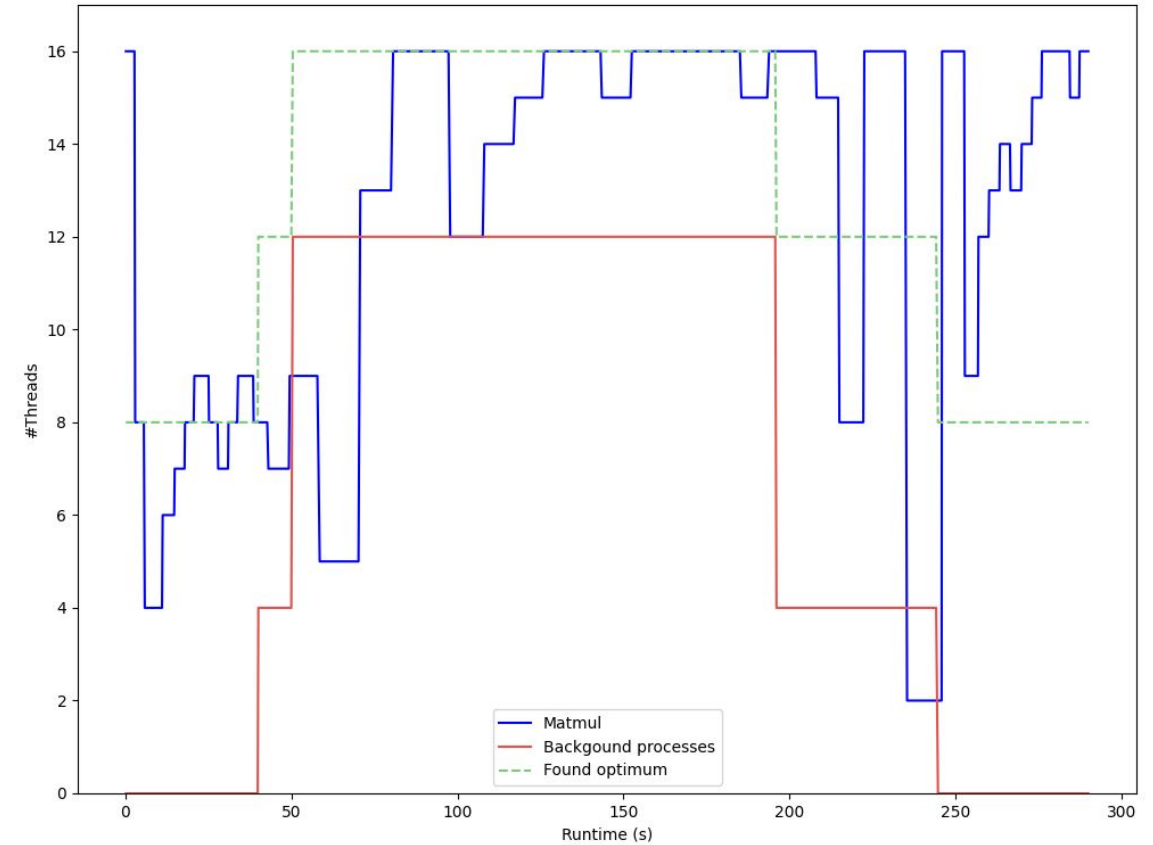


Dynamic vs. Static

# Results

## Energy needed for 1300 matrix multiplications

- Changing background load
- Background load no longer artificial



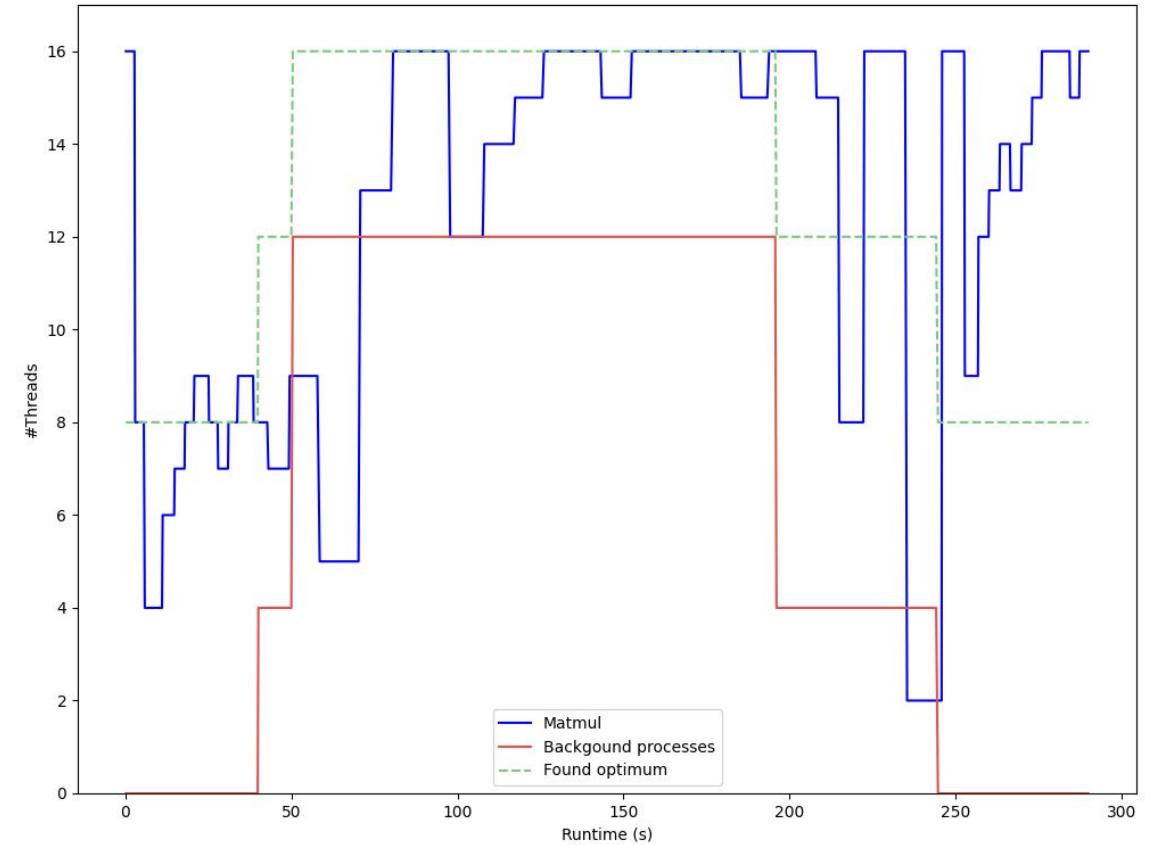
	Dynamic	Static 8 Threads	Static 12 Threads	Static 16 Threads
Energy	12.8kJ	13.6kJ	14.0kJ	14.8kJ
Runtime	290s	302s	303s	322s

Dynamic vs. Static

# Results

## Energy needed for 1300 matrix multiplications

- Changing background load
- Background load no longer artificial



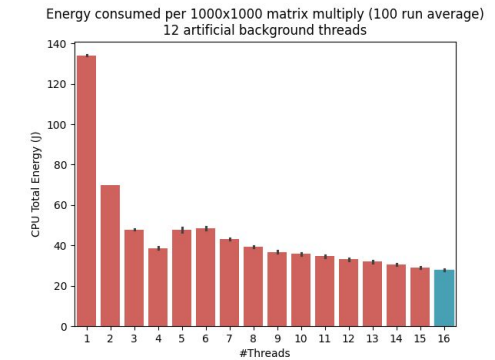
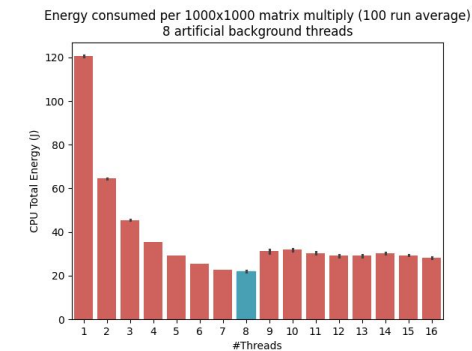
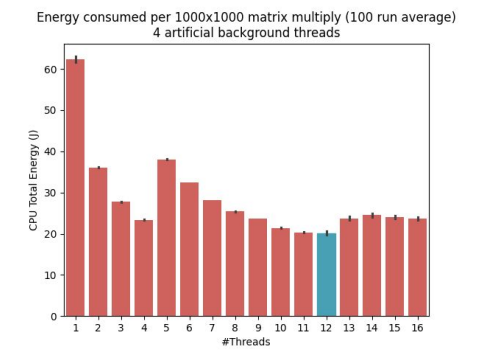
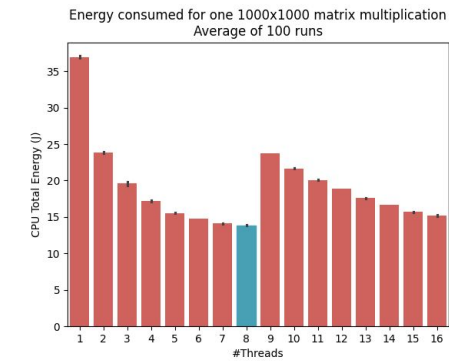
	Dynamic	Static 8 Threads	Static 12 Threads	Static 16 Threads
Energy		-6.3%	-9.4%	-15.6%
Runtime		-4.1%	-4.5%	-11.0%

How much do we lose for being dynamic?

# Overhead

## Magic oracle that knows the optimal thread count

- For the task at hand
- For the given system
- For the given (artificial) background load



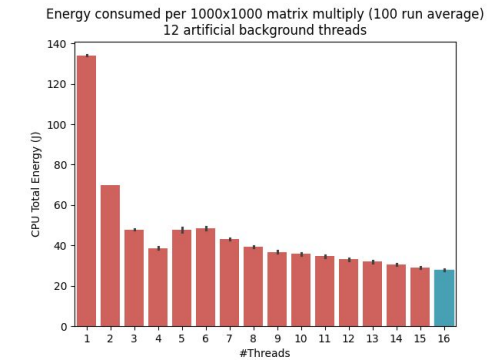
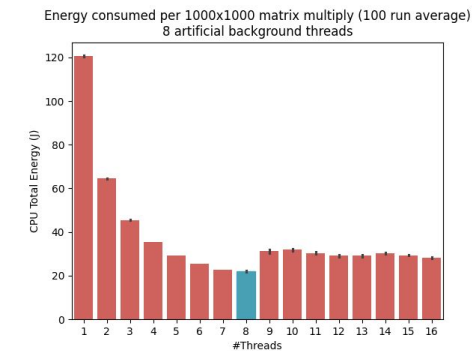
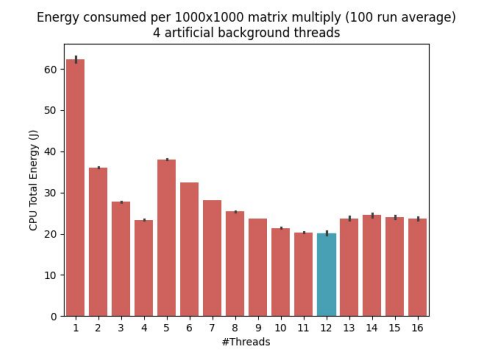
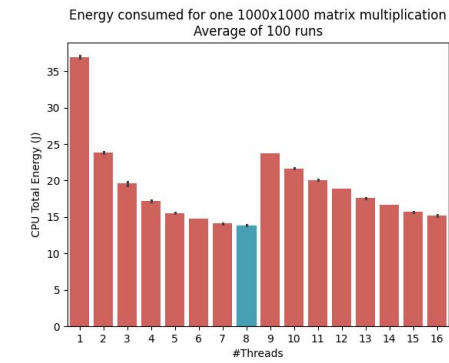
		No BG load	4 BG threads	8 BG threads	12 BG threads
Dynamic	Energy	7.11kj	11.17kj	8.94kj	14.19kj
	Runtime	165s	278s	197s	342s
Oracle	Energy	- *	10.72kj	7.51kj	13.88kj
	Runtime	- *	236s	157s	331s

How much do we lose for being dynamic?

# Overhead

## Magic oracle that knows the optimal thread count

- For the task at hand
- For the given system
- For the given (artificial) background load



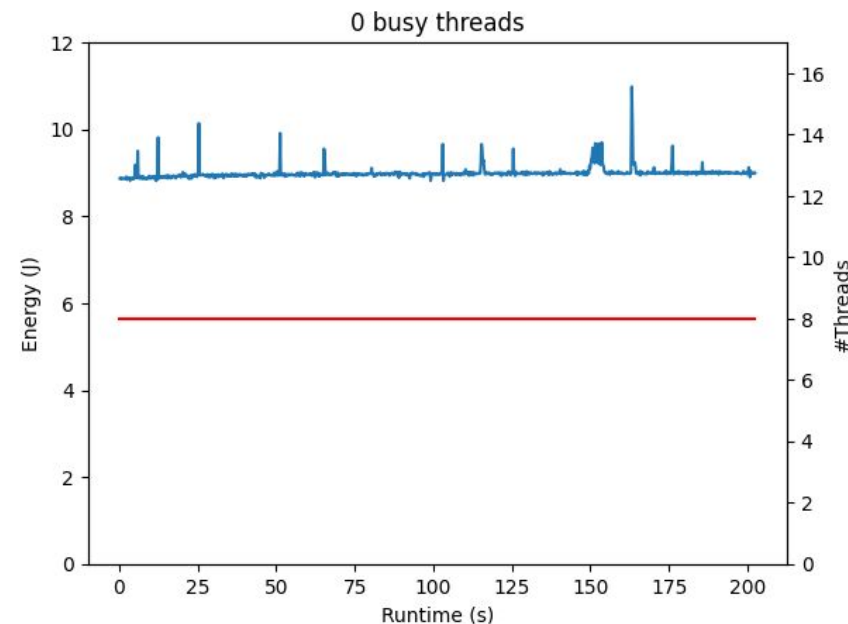
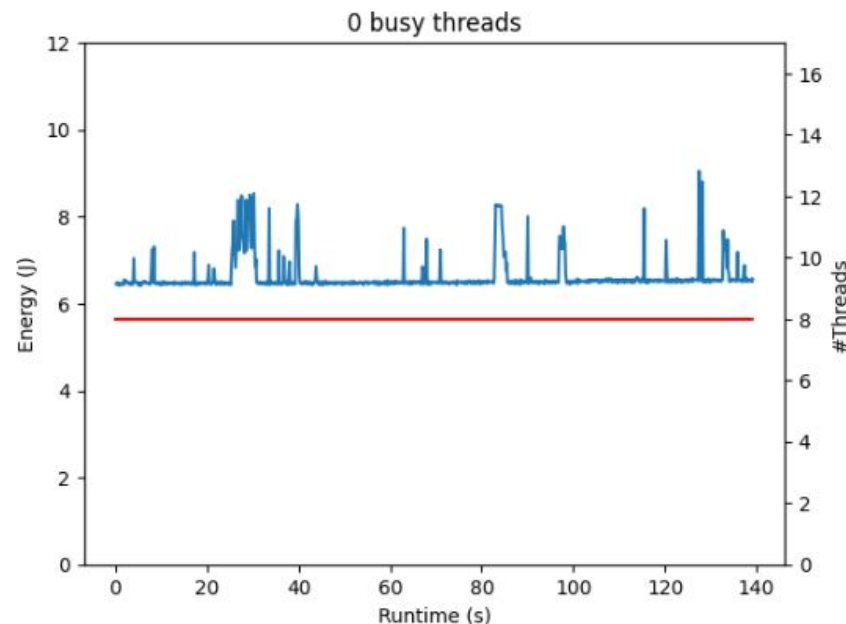
		No BG load	4 BG threads	8 BG threads	12 BG threads
Dynamic	Energy	- *	-4.2%	-19.0%	-2.2%
	Runtime	- *	-17.8%	-25.5%	-3.3%

How much do we lose for being dynamic?

# Overhead \*

**Even on a relatively well controlled device we have a lot of variation**

- Why? Hard to tell!
- Future work: pinning threads might reduce variation
- Anecdotally: less variation with dynamic adaptation



# Conclusion

## Language-independent

- Functional languages especially suitable
- POC in SaC

## Already shows potential to save energy

- Able to adapt to runtime changes
- Reasonable compared to magic oracle
- Other types of programs?
- Other ways to measure?

**6 - 15% less energy consumed vs. static**