

Serverless FPGA

Utilizing dynamic partial reconfiguration in
an FPGA-accelerated FaaS architecture

Martin Lambeck
Advisor: Dr. Atsushi Koshiba
Chair of Computer Systems
<https://dse.in.tum.de/>

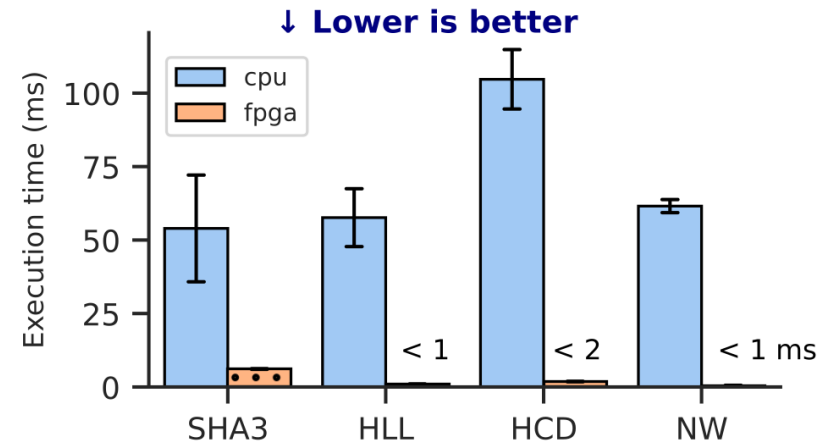


15.07.2023 – 15.07.2024
(part time)

- Motivation
- Design
- Benchmark applications
- Evaluation
- Conclusion

FPGAs

- Offer great computing performance
 - Excel at parallelizable/pipelineable tasks
- Difficult to program and integrate into system



FPGAs

- Offer great computing performance
 - Excel at parallelizable/pipelined tasks
- Difficult to program and integrate into system

Serverless functions (FaaS)

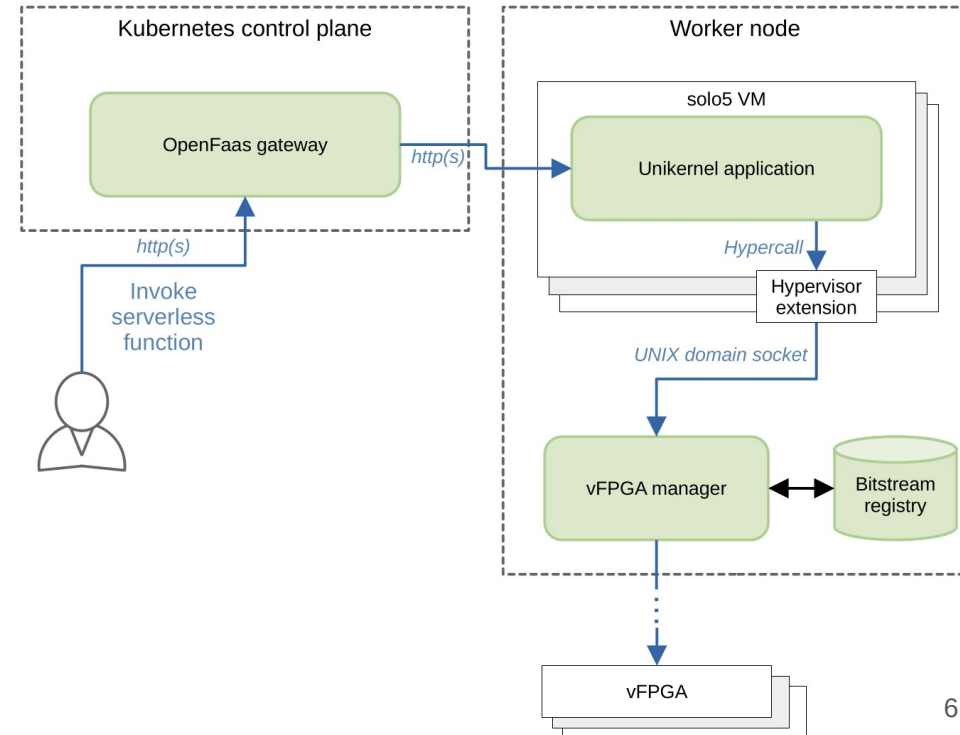
- Popular cloud deployment model
- Simplifies deployment greatly
- No infrastructure management by developer

Can we combine them?

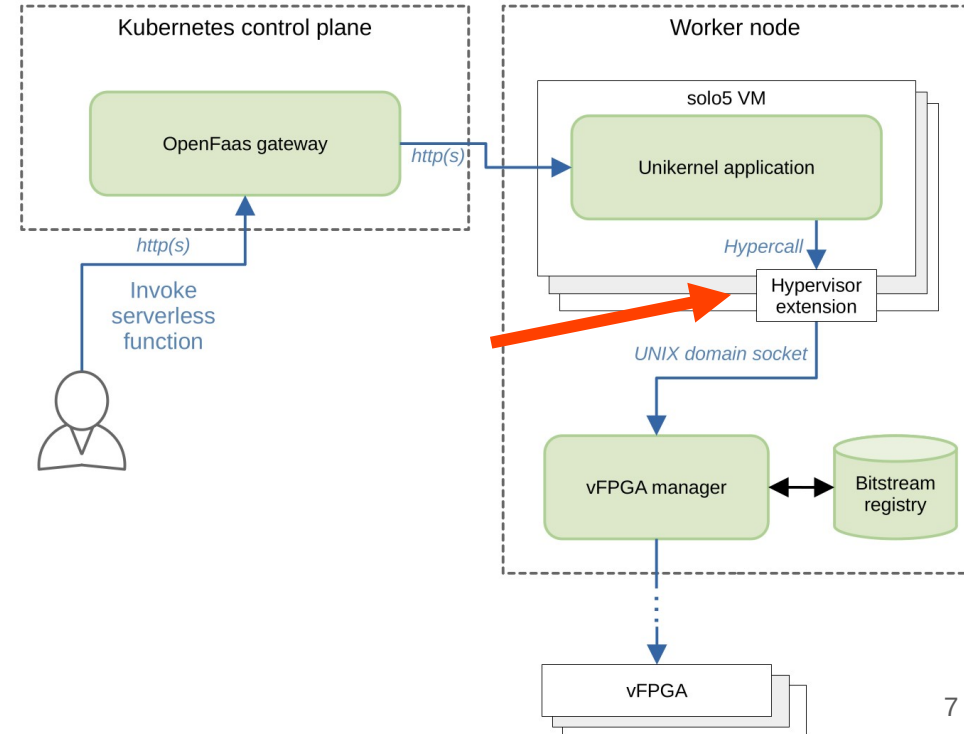
- ~~Motivation~~
- Design
- Benchmark applications
- Evaluation
- Conclusion

Design

- **FPGA:** Use Coyote shell as runtime
 - Multiple isolated slots (vFPGAs)
 - Multi tenancy
- **Host:** A minimal program on the host side is necessary
 - Relays input/output data between FPGA and invoker
 - Unikernel, confined to VM
- **Cluster:** OpenFaaS/Kubernetes for orchestration

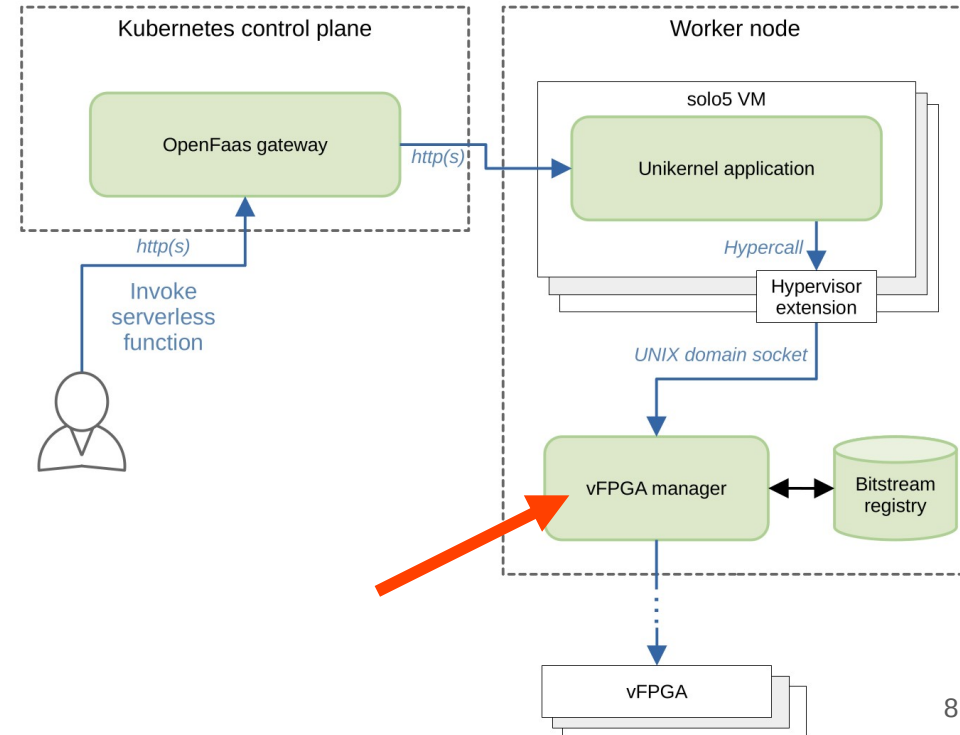


- Hypercalls to exit VM
- ```
solo5_result_t solo5_serverless_set_csr(uint32_t offset, uint64_t value);
solo5_result_t solo5_serverless_load_bitstream(uint32_t config);
solo5_result_t solo5_serverless_map_memory(
 void *addr,
 size_t len,
 size_t input_len,
 size_t output_len);
solo5_result_t solo5_serverless_exec(void);
```
- Communicate with vFPGA manager



# Design

- Client/server architecture
- Receives invocation request from unikernel app
- Schedules invocations
- Manages reconfiguration
- Invoke user logic on FPGA
- Return completion signal





# Outline



- ~~Motivation~~
- ~~Design~~
- Benchmark applications
- Evaluation
- Conclusion

- Select popular algorithms
  - Port/implement FPGA implementation to Coyote
  - Baseline: CPU implementation
- 12 benchmark applications:

|                  |                         |
|------------------|-------------------------|
| AddMul           | AES (ECB mode)          |
| SHA256           | SHA3                    |
| GZIP             | MATMUL (64x64)          |
| Needleman-Wunsch | hls4ml                  |
| Hyperloglog      | Harris Corner Detection |
| MD5 brute force  | FFT Auto-Correlation    |

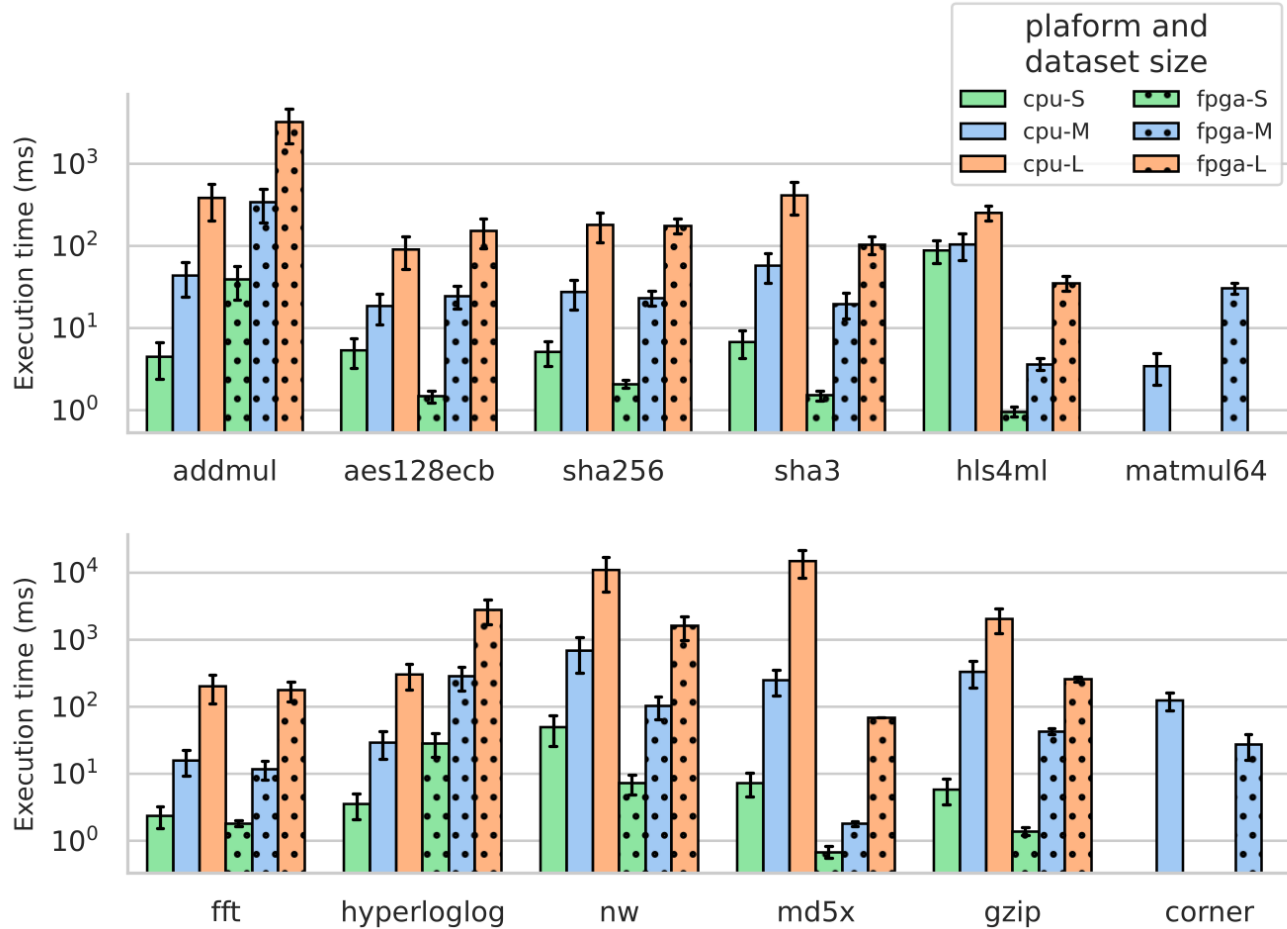
# Outline



- ~~Motivation~~
- ~~Design~~
- ~~Benchmark applications~~
- Evaluation
- Conclusion

- End-to-end evaluation
  - Evaluate all benchmarks
  - Measure duration from moment input data has been received until result is ready
  - No network overhead!
- Micro-benchmarks
  - Measure impact of huge pages
  - Measure parallel efficiency
  - Measure reconfiguration overhead

# Results (End-to-end)

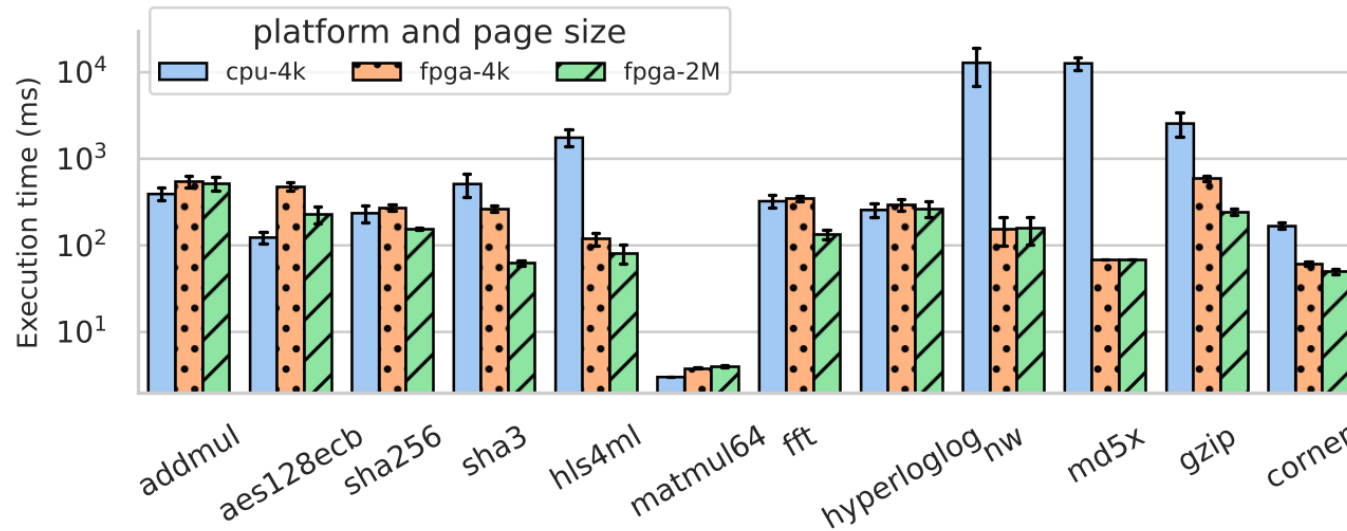


## FPGA speed-up factors

|                 | S    | M    | L    |
|-----------------|------|------|------|
| <b>addmul</b>   | 0.12 | 0.13 | 0.12 |
| <b>aes128</b>   | 3.7  | 0.76 | 0.6  |
| <b>sha256</b>   | 2.4  | 1.2  | 1.0  |
| <b>sha3</b>     | 4.5  | 2.9  | 4.0  |
| <b>hls4ml</b>   | 91   | 29   | 7.2  |
| <b>matmul64</b> |      | 0.11 |      |
| <b>fft</b>      | 1.3  | 1.4  | 1.1  |
| <b>hll</b>      | 0.12 | 0.10 | 0.11 |
| <b>nw</b>       | 6.9  | 6.7  | 6.9  |
| <b>md5x</b>     | 11   | 140  | 220  |
| <b>gzip</b>     | 4.2  | 7.8  | 8.0  |
| <b>corner</b>   |      | 4.6  |      |

# Results (Huge pages)

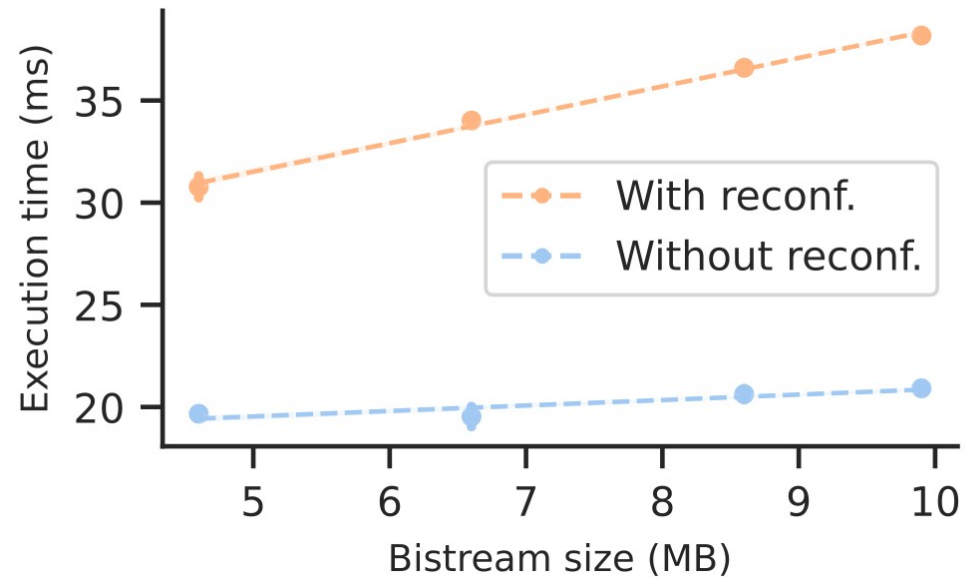
FPGA + huge pages  
speed-up factors



| Baseline:       | cpu-4k | fpga-4k |
|-----------------|--------|---------|
| <b>addmul</b>   | 0.76   | 1.1     |
| <b>aes128</b>   | 0.54   | 2.1     |
| <b>sha256</b>   | 1.5    | 1.7     |
| <b>sha3</b>     | 8.3    | 4.2     |
| <b>hls4ml</b>   | 22     | 1.5     |
| <b>matmul64</b> | 0.76   | 0.95    |
| <b>fft</b>      | 2.4    | 2.6     |
| <b>hll</b>      | 0.97   | 1.1     |
| <b>nw</b>       | 82     | 0.98    |
| <b>md5x</b>     | 180    | 1.0     |
| <b>gzip</b>     | 11     | 2.4     |
| <b>corner</b>   | 3.4    | 1.2     |

# Results (Reconfiguration)

- Dynamic partial reconfiguration induces overhead
- Linear correlation between bitstream size and duration



| Clock regions  | 8      | 16     | 24     | 32     |
|----------------|--------|--------|--------|--------|
| Bitstream size | 4.6 MB | 6.6 MB | 8.6 MB | 9.9 MB |

# Outline



- ~~Motivation~~
- ~~Design~~
- ~~Benchmark applications~~
- ~~Evaluation~~
- Conclusion



# Conclusion

- FPGAs offer excellent performance for some workflows
- Coyote shell enables multi-tenancy on the FPGA
- FaaS simplifies infrastructure management
- => FaaS + FPGA is a feasible deployment model

Thank you!  
Questions?

