# ETH *zürich*

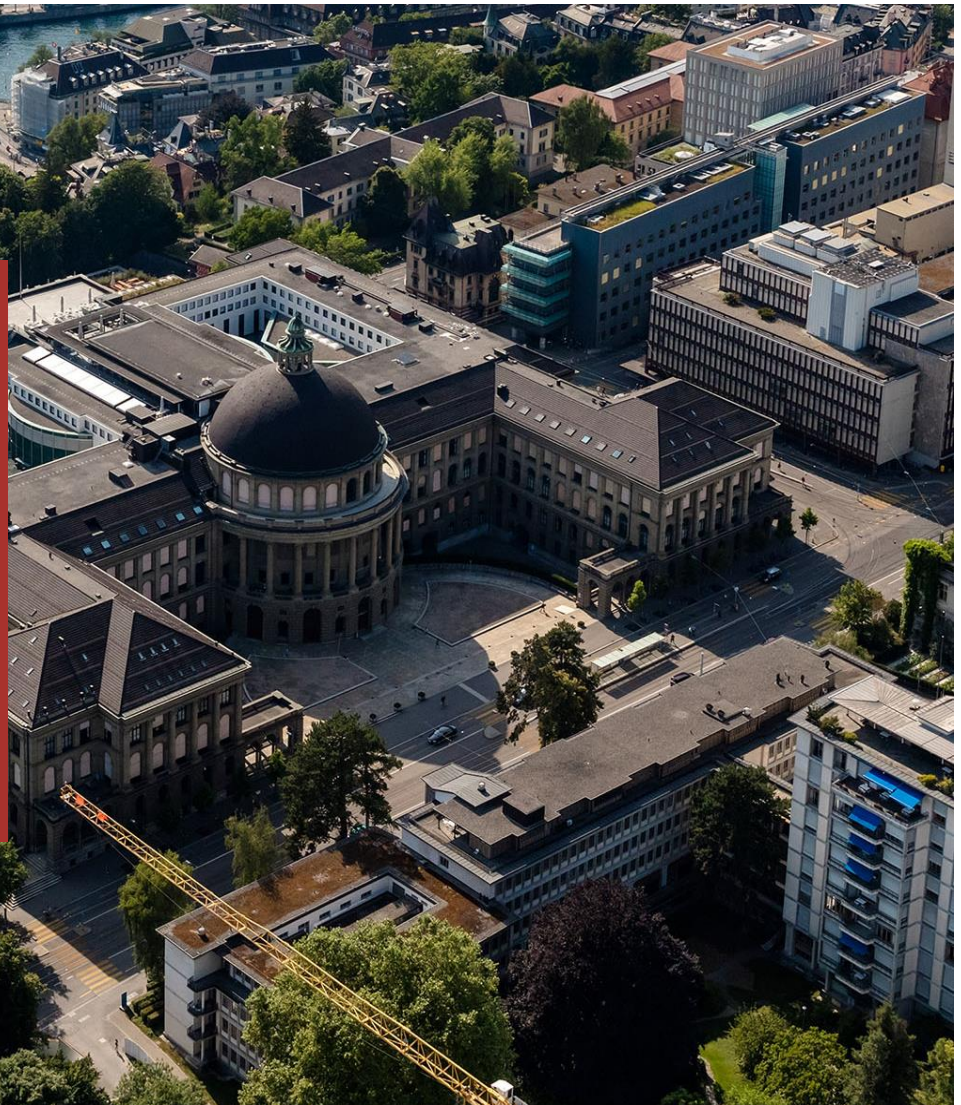# Networking Abstractions for Modern
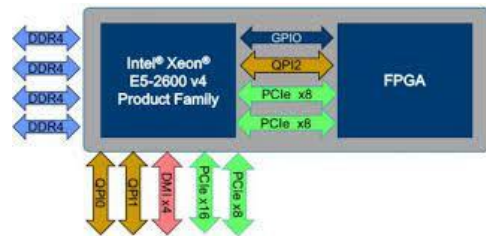# Heterogeneous Systems

Dario Korolija
*Systems Group, Dept. of Computer Science, ETH Zurich*
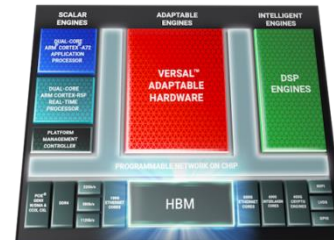
**D** INFK

*Systems@ETH Zürich*
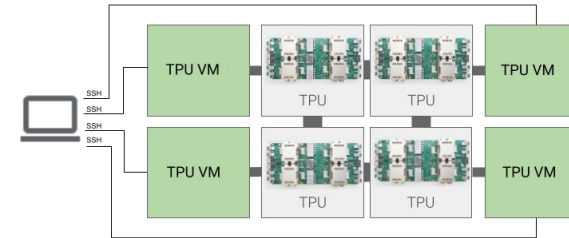
# Introduction...

❖ **Scpecialized hardware becoming a reality**

➢ *Amazon*, *Microsoft*, *Google*, *Alibaba*, *Intel*, *AMD* ...
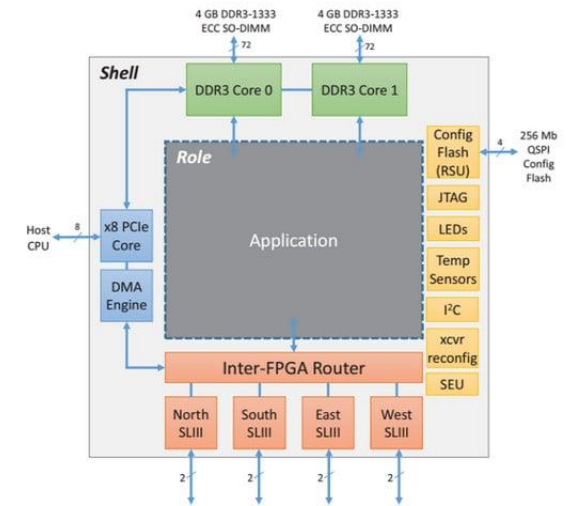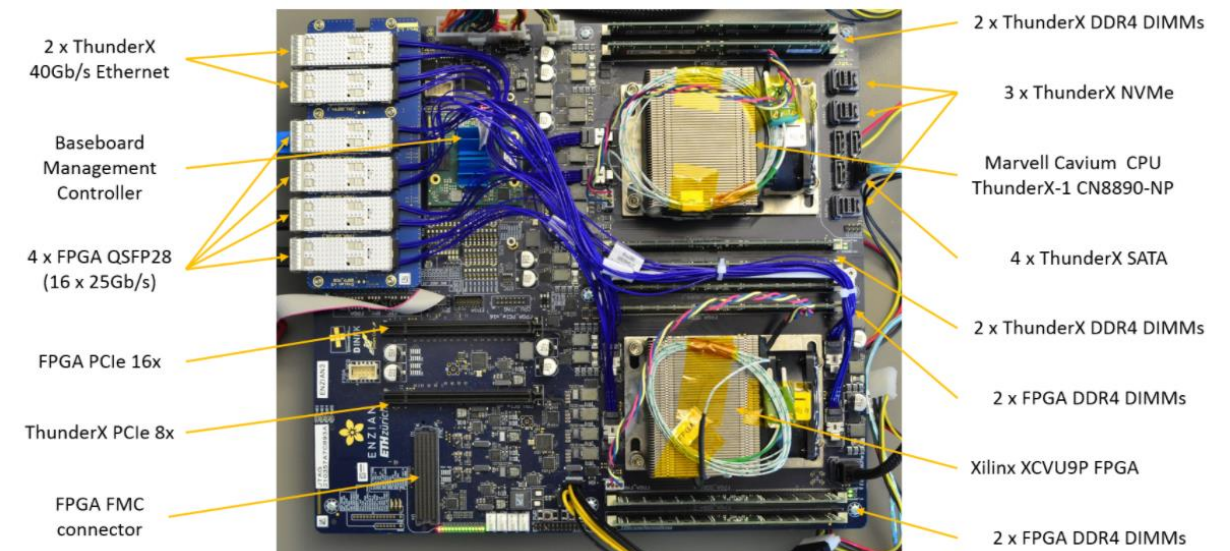


Intel HARP



AMD Versal



Google TPU



Microsoft Catapult
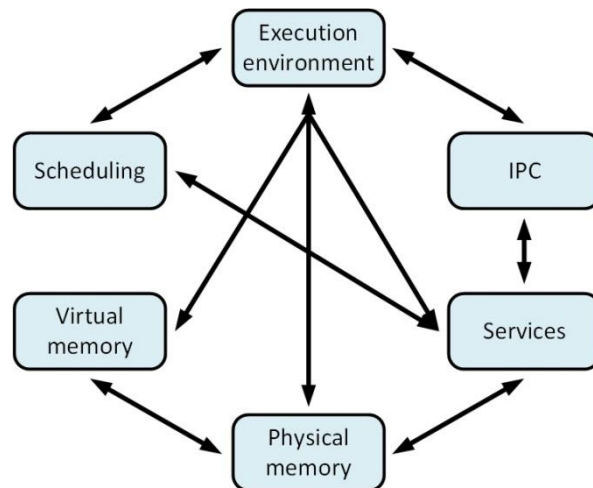
❖ **One built within System group at ETH:**



ENZIAN

# Microkernel for FPGAs
## Hybrid computing system

❖ Plenty of research, focused on individual functionalities only

❖ Coyote provides a complete minimal core set of essential features on which further services can be used

Interdependence

# Coyote

## System Architecture

❖ Microkernel for FPGAs

➢ Multiple isolated untrusted vFPGAs

➢ Spatial and temporal sharing (PR)

➢ Dynamic reconfiguration (scheduler)

➢ Virtualized memory

➢ Unified memory (HMM)

➢ HBM and DRAM striping service

➢ Shared TCP/IP service

➢ Shared RDMA service

➢ Unified logic interface –> portability

➢ RTL and HLS support

➢ Runs on u50, u200, u250, u280, u55c, vcu118, *Enzian*



```
/**
 * No unified memory example
 */
void *data, *data_fpga;

data = malloc(size);
data_fpga = fpga_malloc(size);

prepare_data(data, size);

offload_data_to_fpga(data, data_fpga, size);

execute_fpga_kernel(data_fpga, size);

sync_data_from_fpga(data, data_fpga, size);

print_results(data);

fpga_free(data_fpga);
free(data);
```

```
/**
 * Unified memory example
 */
void *data;

data = malloc(size);

prepare_data(data, size);

execute_fpga_kernel(data, size);

print_results(data);

free(data);
```

# Coyote

TCP/IP stack

❖ Open source TCP/IP stack

❖ Added support for shared functionality between all applications running within an FPGA

# Coyote

RDMA stack

❖ Open source RDMA stack (UC, RC)

❖ RDMA over Converged Ethernet (RoCE v2)

❖ Implemented on top of UDP/IPv4/IPv6 (far lower overhead than iWARP)

❖ InfiniBand (IB) transport packets over Ethernet (READ, WRITE, SEND)



https://docs.nvidia.com/networking/display/WINOFv55053000/RoCEv2

# Coyote
RDMA advantages

- ❖ Bypasses kernel space
- ❖ Zero-copy data movement
- ❖ Cheap pipelined processing (directly on the NIC)



ETH zürich    D INFK

# Coyote

SW

❖ Layered parallelization potential

❖ User space abstractions

➤ <u>cSched</u> – Coyote scheduler, reconfiguration controller

➤ <u>cProc</u> – Coyote process, multiple can run within a single *vFPGA*

➤ <u>cThread</u> – Coyote thread, multiple can run within a single *cProc*. Task level parallelism

➤ <u>cTask</u> – Coyote task, arbitrary user variadic function executed by *cThreads*

➤ <u>cService</u> - Coyote library daemon, background service, UDS for IPC

```
.
.
.
/**
 * Open a Unix Domain Socket and send a decrypt_and_compress task request
 * This is the only place of interaction with Coyote
 */
cLib clib("/tmp/coyote-daemon-vfid-0);
clib.task({opDecryptAndCompress, {mem, size, key}}); // blocking
.
.
.
```

**ETH** *zürich*    **D** INFK

# Coyote
Repository

❖ Github: https://github.com/fpgasystems/Coyote.git

❖ Internal and external users

❖ Example designs (perf. runs, rdma, tcp/ip, rpc, hbm, dram, services)

❖ Documentation …

ETH*zürich*    **D** INFK

# Coyote
## Current Work

❖ Virtualization – can we pull vFPGAs all the way to VMs layer?

    ❖ Virtual Function I/O (VFIO) Mediated devices

# Coyote

Current work

- ❖ Nested reconfiguration (*Nested-DFX*)
  - ➢ <u>Static</u> layer
  - ➢ <u>Service</u> layer
  - ➢ <u>Dynamic</u> layer
  - ➢ <u>Application</u> layer

ETH *zürich*    **D** INFK

# Advantages of the system

❖ Networking advantages:
  ❖ Streaming interfaces:
  ❖ Kernel invocation overhead XRT~50us, Coyote: ~1-1.5us
  ❖ RDMA

# Farview
## Motivation

❖  I/O overheads – main bottleneck

❖  More and more data in local DRAM

❖  Excessive data movement

❖  Memory capacity limitations

Disaggregation of
of computer and storage and storage



Compute Nodes

Smart Memory Nodes

ETH zürich    D INFK

# Farview

Introduction

- ❖ FPGA-based smart NIC making DRAM available as a pool of network attached memory accessible on demand over high performance RDMA network.

- ❖ Performs line-rate query processing with minimal overheads

- ❖ Farview is a disaggregated buffer cache with operator pushdown capabilities

# Farview
## Overview

❖ Farview addresses inefficient data movement and memory capacity limitations

❖ Consider the following queries:

SELECT T.a, S.b
FROM T, S
WHERE T.id = S.id
AND T.c > 50 AND S.d < 2012;

SELECT R.d, S.b
FROM R, S
WHERE R.id = S.id
AND R.a = 3.14 AND S.a <> 2020;

❖ Farview centralizes the buffer cache and performs operator pushdown

ETH *zürich*      D INFK

# Farview

Example

❖ processing in Farview, simple RDMA READ operation:



❖ AES decryption on the same data as it is being read:

# Farview
## Example



Throughput [GBps] vs Transfer size [bytes]

Legend:
- RDMA READ
- Farview READ + AES Decryption

ETH zürich    D INFK

# Farview

## Architecture

❖ Several components needed:
  ➢ DRAM (HBM)
  ➢ Memory controllers
  ➢ Memory management unit
  ➢ Network stack
  ➢ Mechaniscm for concurrent access
  ➢ Stream processing capacity
  ➢ Mechanism to swap operators

❖ Three distinct layers
  ➢ Operator stack
  ➢ Memory stack
  ➢ Network stack

# Farview

## Operator stack



➢ Operator stack split into multiple isolated dynamic regions that operate concurrently

➢ Operator pipeline can execute a set of queries

**A single dynamic region and interfaces:**

# Farview
## Operators



- ❖ Operator pipeline
- ❖ Farview currently supports a range of operators (row store):
  - ➢ Projection operators (smart addressing, projection)
  - ➢ Selection operators (selection, regex, vectorized selection)
  - ➢ Grouping operators (distinct, group by)
  - ➢ System operators (encryption/decryption, parsing, packing)

**ETH** zürich   **D** INFK

# Farview

## Memory stack

➢ Implements the actual memory buffer pool

➢ Organized into multiple channels

➢ Interleaving abstraction to aggregate the bandwidth

➢ Can process data at higher rates than the available network bandwidth



Memory stack architecture:

ETH zürich     Systems@ETH Zürich     D INFK

# Farview
## Network stack

- ➢ Manages all external connections and requests for all concurrent accesses
  - ➢ Supports RoCE v2 at 100Gbps
  - ➢ Open source network stack[2]

- ➢ Special Farview verb based on InfiniBand SEND for query requests

- ➢ Comparable latencies to one-sided RDMA verbs

[2] **StRoM**: *Smart Remote Memory, EuroSys '20*
David Sidler, Zeke Wang, Monica Chiosa, Amit Kulkarni, Gustavo Alonso

**Network stack architecture:**

# Farview

Programmatic Interface

- High level data API covering both the
  - critical path operations and
  - connection management operations

- Written in C++

- Intended to be used by Farview query compiler.

- • bool openConnection(Qpair *qp, Fview *node);

- • bool loadPipeline(Qpair *qp, int32_t opid);

- • void tableRead(Qpair *qp, Ftable *ft);

- • void tableWrite(Qpair *qp, Ftable *ft);

- …

- • void farView(Qpair *qp, Ftable *ft, uint64_t *params);

- • void select(Qpair *qp, Ftable *ft, uint64_t *proj_flags, uint64_t *sel_flags, float predicate) {

      …

          farView(qp, ft, params);

  }

# Farview
## Implementation

- Farview supports a range of FPGA data center cards (Alveo u50, u55c, u200, *u250*, u280, Enzian)
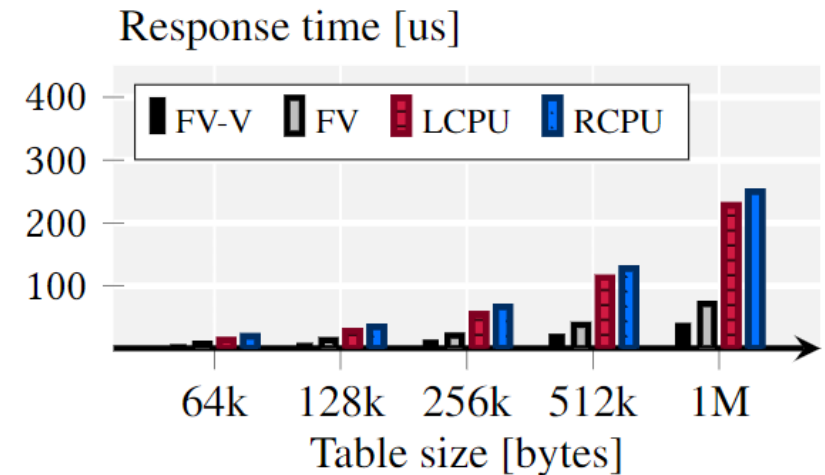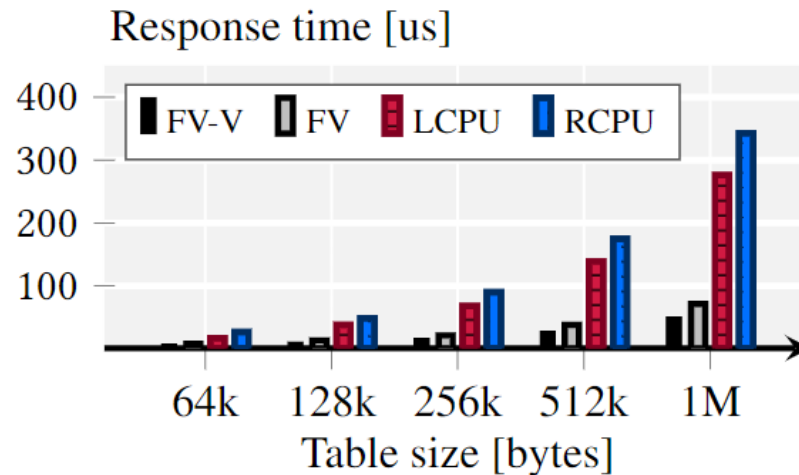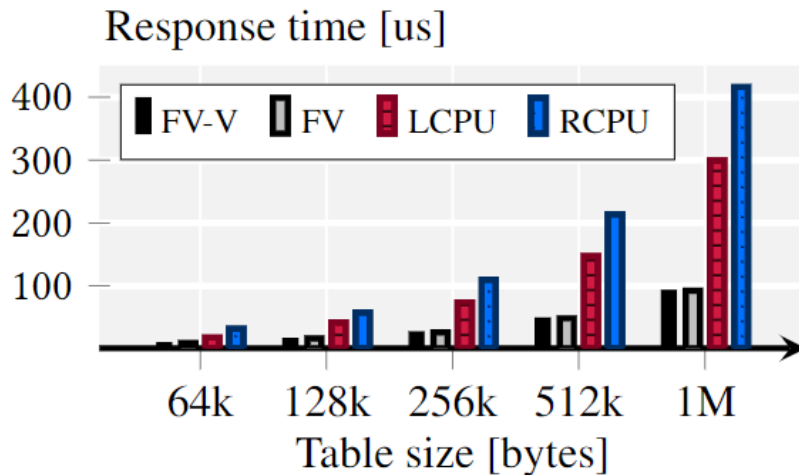- Low resource usage:

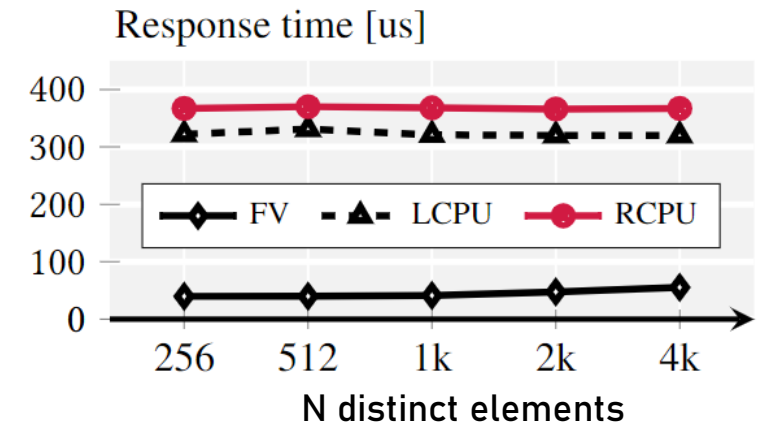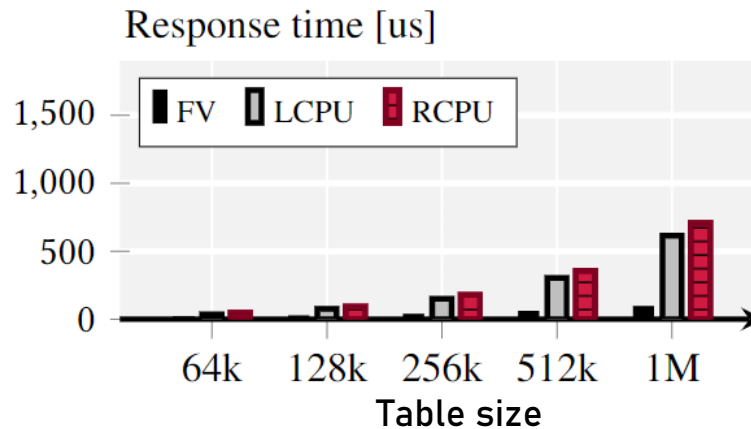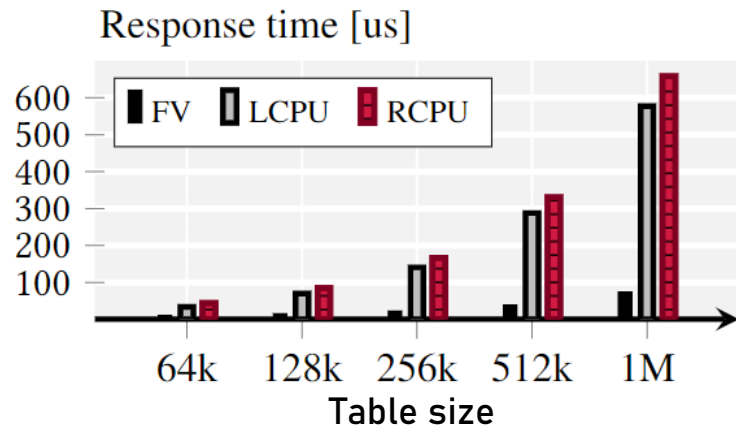| Configuration | CLB LUTs | Regs | BRAM tiles | DSPs |
|---|---|---|---|---|
| 6 regions | 24% | 23% | 29% | 0% |
| **Operators  (per dynamic region)** | **CLB LUTs** | **Regs** | **BRAM tiles** | **DSPs** |
| Projection / Selection | < 1% | < 1 % | 0% | 0% |
| Regex engine | 2.3 % | < 1% | 0% | 0% |
| Distinct / Group by | 2.1% | 1.3% | 8% | 0% |
| En(de)cryption | 3.6% | < 1% | 0% | 0% |
| Packing / Sending | < 1% | < 1% | 0% | 0% |

ETH zürich          D INFK

# Farview

## Benchmarks

❖ Selection performed across different selectivity levels

❖ Farview outperforms two baselines (traditional database, remote memory) across different selectivity levels

Response times for selection queries with 100%, 50% and 25 % selectivity, respectively:
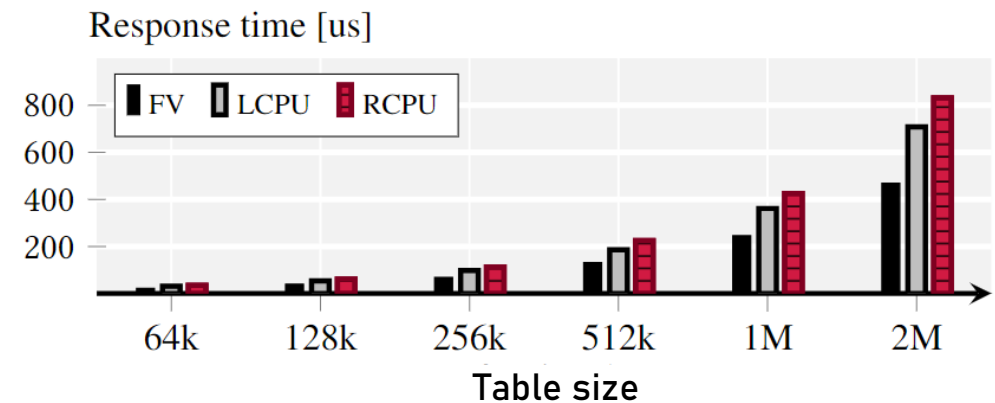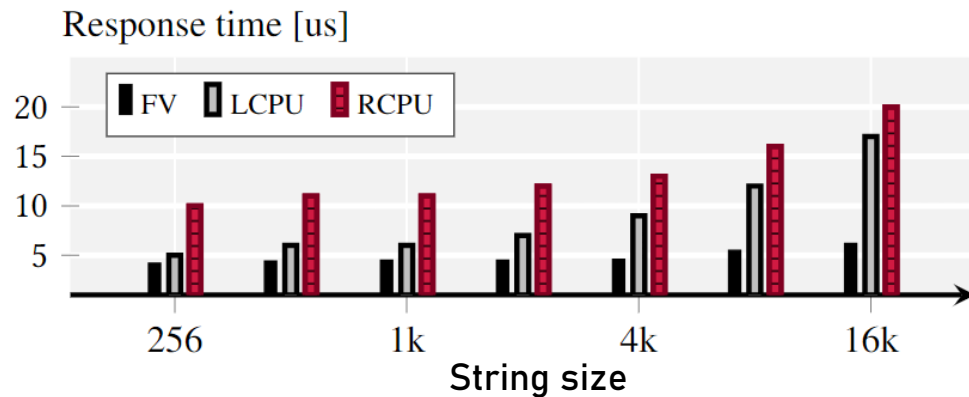
# Farview

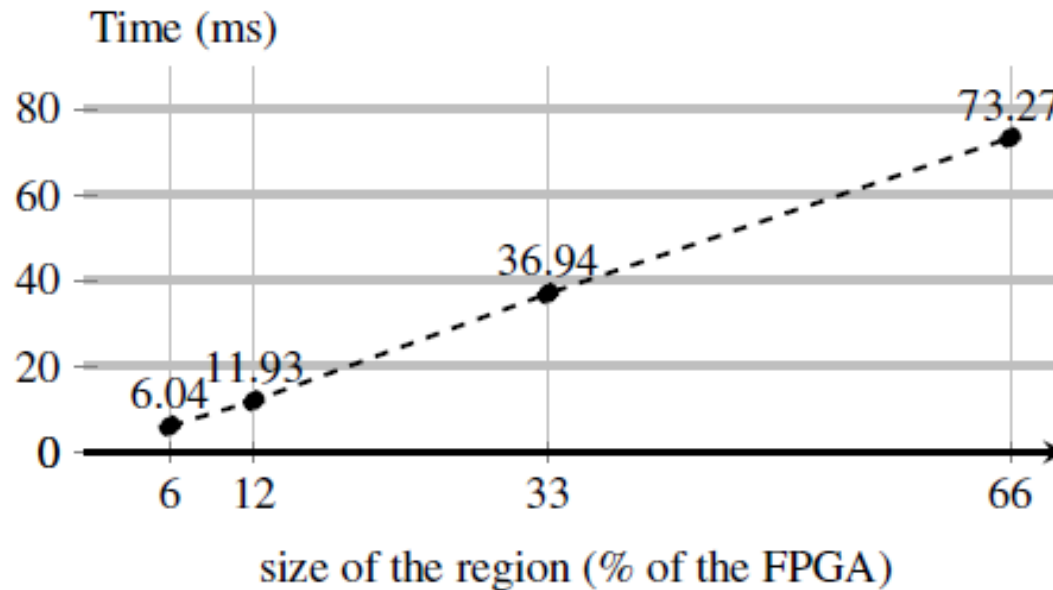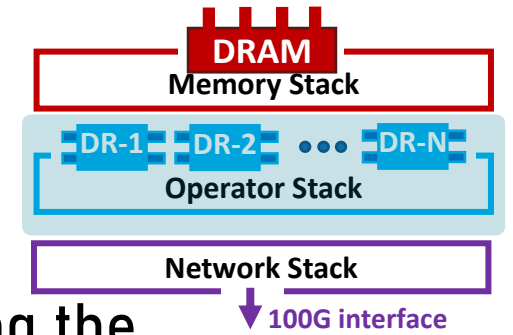❖ **Distinct, Group by, Group by (sweep by n. of distinct el.)**



❖ **Regex and concurrent queries (3x)**

# Farview
## Operator Pipeline Swap

DRAM
Memory Stack

DR-1 DR-2 ••• DR-N
Operator Stack

Network Stack

100G interface

- Operator pipelines can be swapped during *runtime* without affecting the integrity of the system
- Gives Farview a much needed flexibility in comparison to traditional accelerators
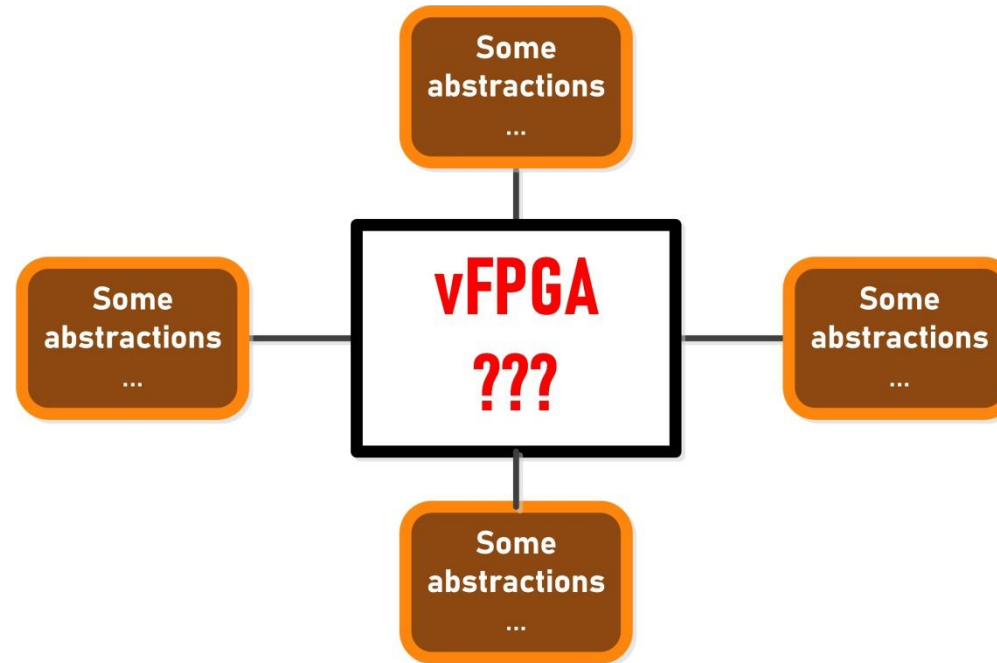- Swap time in the order of milliseconds:



ETH zürich   Systems@ETH zürich   D INFK

# Farview
## Current work

❖ Additional operators …

❖ Interaction with the storage layer (deduplication …)

❖ Scale out …

❖ Serverless (end-to-end latencies in s range)

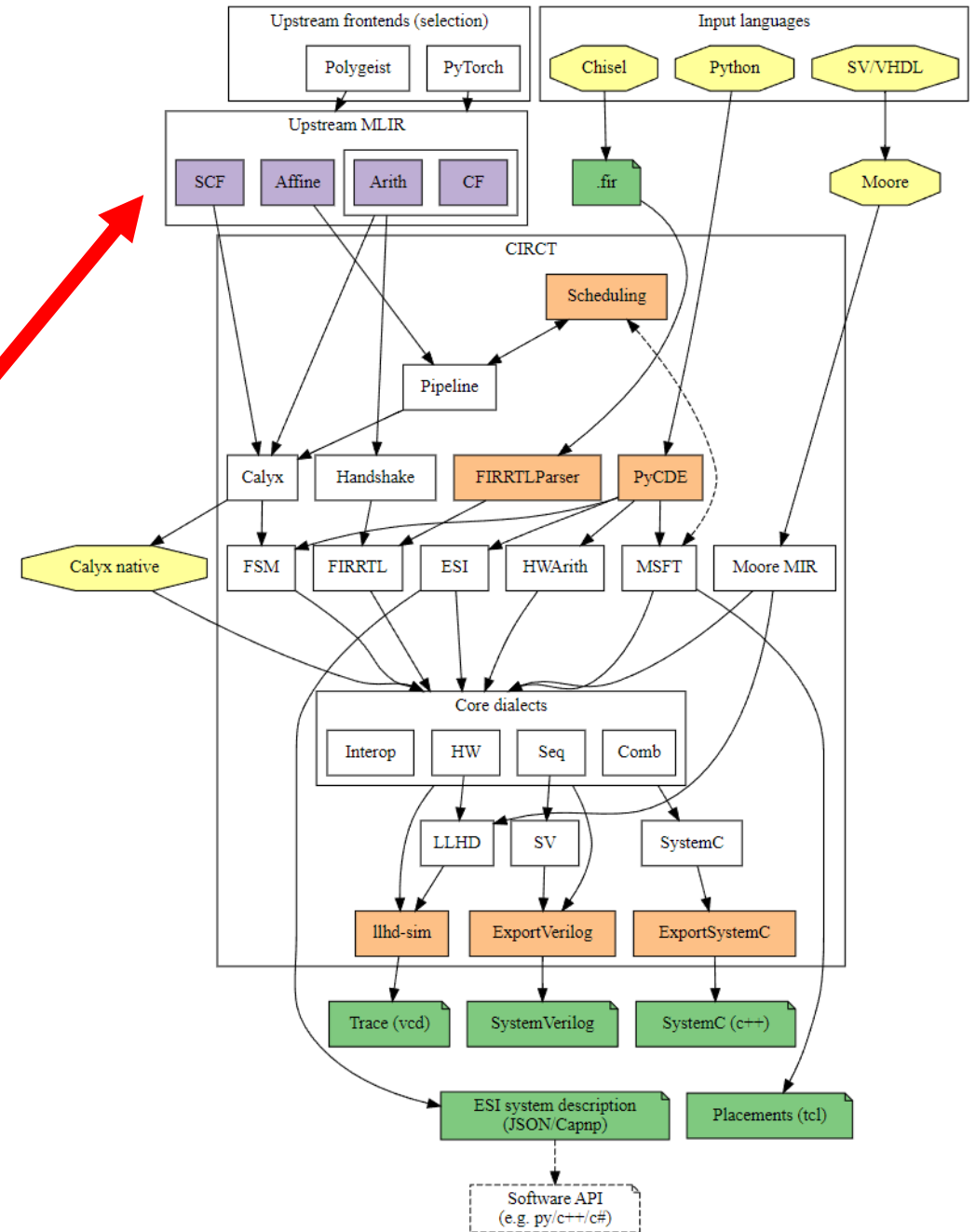❖ Proper database frontend (Modularis)

**ETH**zürich          **D** INFK

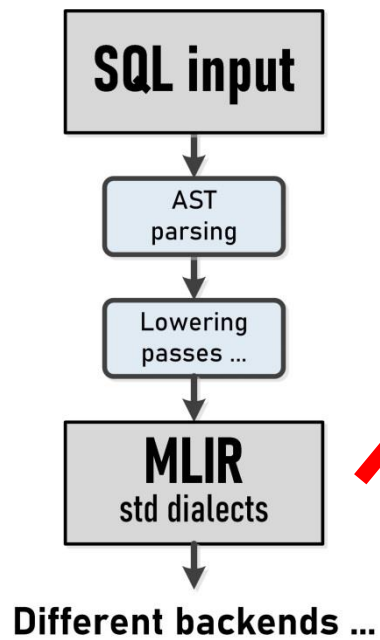# What about the vFPGA?

❖   HDL … (SVerilog, VHDL, HLS, OpenCL ,…)

# MLIR compiler

- MLIR is a novel approach to building a compiler

- CIRCT built in MLIR, targets HDL
  (https://github.com/llvm/circt.git)

- Modularis infrastructure being ported to MLIR

# Modularis + Farview

- 2 Master thesis projects
- Stream-dialect => Coyote-CIRCT (*https://github.com/fpgasystems/Coyote-CIRCT.git*)

# Questions?

ETH zürich          D INFK