

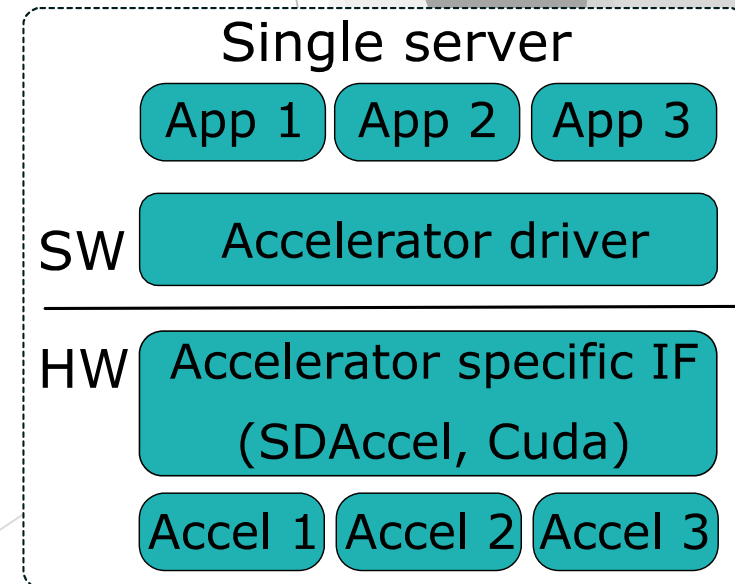


# Simplifying Software Access and Sharing of FPGAs in Datacenters

Stelios Mavridis, **Manos Pavlidakis**, Ioannis Stamoulias, Christos Kozanitis, Nikos Chrysos,  
Christoforos Kachris, Dimitrios Soudris, and Angelos Billas

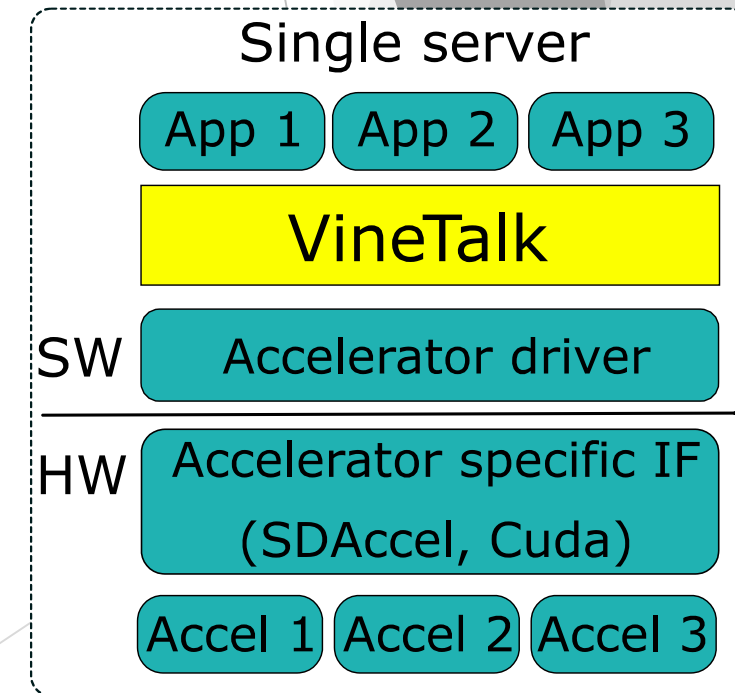
# Challenges for using accelerators in servers

1. Difficult to write Host code for single & different accel. type(s)
  - Accelerator code and Host code are currently tightly coupled
    - Results in increased programming effort for application developer
    - Reduced accelerator code reuse
    - Requires detailed programming knowledge for Host & Accelerator side
  
2. Difficult to share the same IP core from multiple apps
  - Typically IP cores are dedicated
  - But in consolidated servers sharing is important
  
3. Difficult to use multiple accelerators from one app
  - Host code has to change significantly



# VineTalk: Software layer between FPGAs & apps

- VineTalk addresses the pre mentioned issues
  - Host code and FPGA code is decoupled
    - Host code is written once, regardless of the accelerator number & type
  - Provides IP core sharing from multiple applications
  - Apps can use multiple/heterogeneous accelerators
- VineTalk supports VMs, Native, and Containers
- VineTalk consists of two main components
  - Transport protocol
  - Extensive scheduling
- VineTalk *virtualizes* accelerators



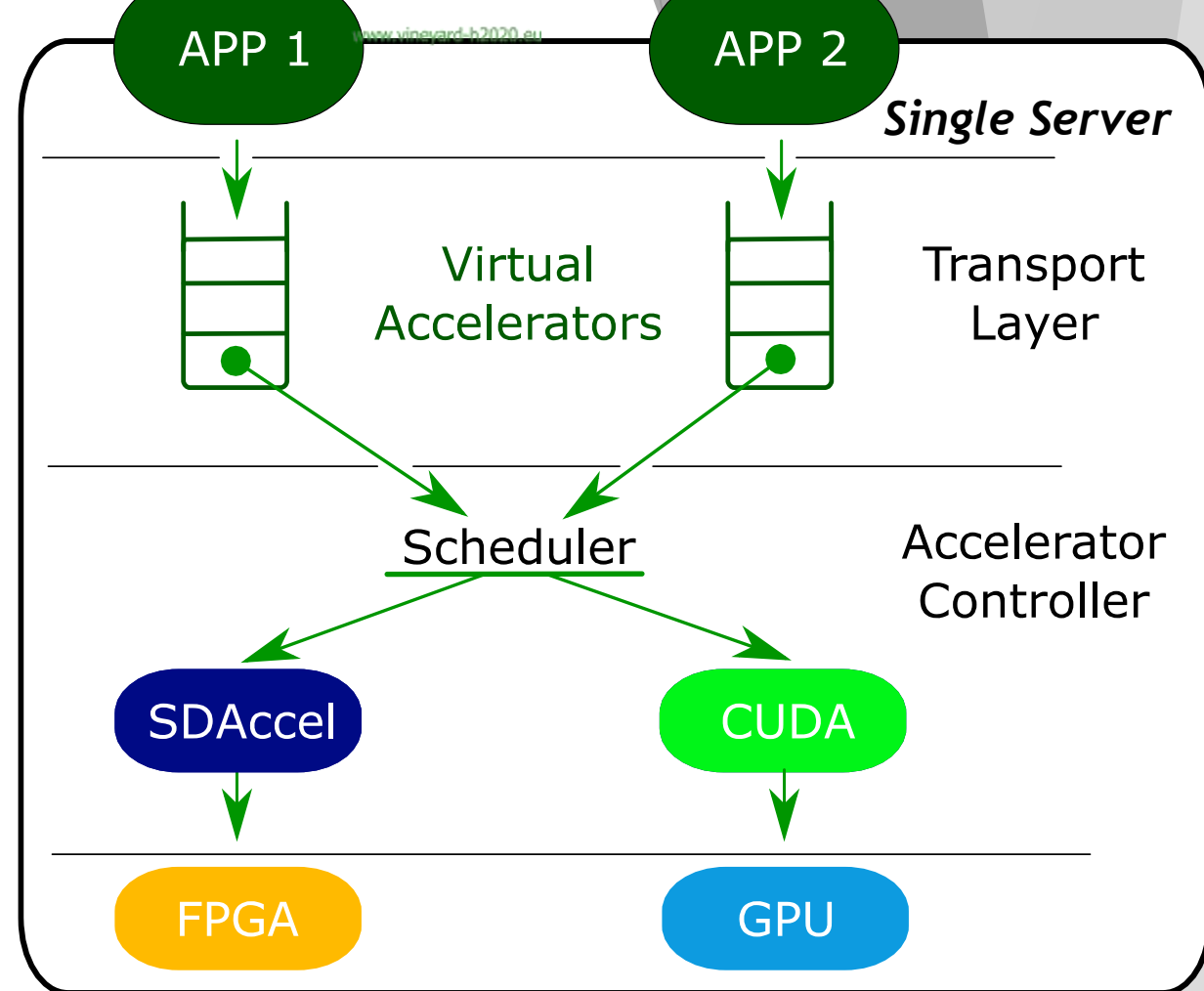
# VineTalk's design

## ○ Transport layer

- Implemented as shared memory
  - Enables VMs & Containers
  - Faster than network approaches
- Virtual Accelerators
  - Implemented as task queues
  - Allows FPGA sharing

## ○ Accelerator Controller

- Schedules **multiple** apps over a **single** accelerator
- Schedules **one** app over **multiple** accelerators
- Schedules **one** app over **heterogeneous** accelerators



# Preliminary evaluation

- Perform risk analysis
  - With three financial apps with & without VineTalk
  - Black&Scholes, Black-76, and Binomial
- Question 1: Is VineTalk expensive ?
  - Up to 4% slower compared to native execution
- Question 2: Is FPGA sharing expensive ?
  - With 2 concurrent apps
  - 2% less task rate compared to 1 app running standalone
- Programming effort
  - Decrease the lines of **Host** code up to 30% compared to native
  - Accelerator specific code moved to Accelerator Controller

# Summary

- VineTalk virtualizes heterogeneous accels. in consolidated servers
  - Using one efficient transport layer and an Accelerator Controller
  - The Accelerator Controller provides task scheduling and FPGA sharing
- Benefits of VineTalk
  1. Accelerator sharing
  2. Host code is written once
  3. Apps can be executed in a Single/Multiple/Heterogeneous accelerators
  4. Apps can run in VMs, Natively, and in Containers
- Our preliminary results show low overhead in simple case

▶ Thank you! Questions?