



Full Circuit FHE Acceleration on Cornami Fabric

Dr. Vineet Chadha

FHE Engineering, Cornami

- The FHE Problem Statement
- Cornami FHE Solution
- Cornami Architecture Overview
- Secure Computing Interface Framework (SCIFR API™)
- Summary and Future Directions

- Fully Homomorphic Encryption (FHE) is de-facto solution to secure the in-use data in an untrusted environment
 - Data, once encrypted, need not be decrypted for processing
 - Data can be decrypted with a private key in a trusted environment for viewing
- **Problem Statement:** The challenge of FHE is that it requires 10^6 greater performance than x86 and conventional processors
 - Performance gap between encrypted data vs unencrypted data processing



Fully Homomorphic Encryption

- Today's monolithic processors have bounded limits and **do not scale** in computational performance and/or I/O bandwidth. They also are inefficient in power use.
- A typical FHE application would require racks of computers in datacenter with the workload running for days



The Fully Homomorphic Encryption Solution

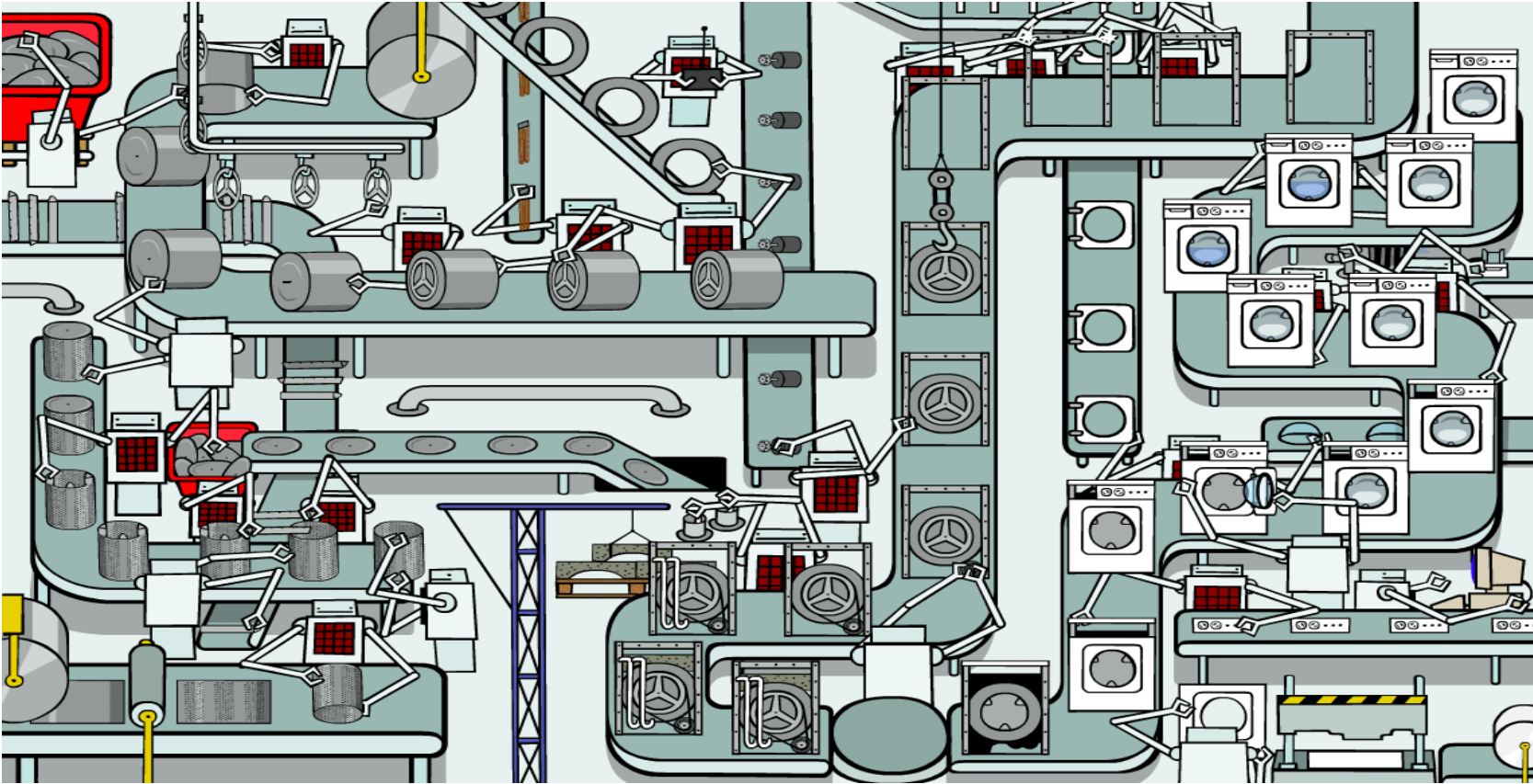
- Requires an execution platform that can:
 - Seamlessly address the multiple security guarantees for an FHE Application
 - Scale up or scale down system resources such as computational performance, I/O bandwidth, and storage
 - Scale up or scale down FHE configuration parameters
 - Process arbitrary precision arithmetic as required by FHE applications
 - Configure resources on a per application basis **on-demand**



The Fully Homomorphic Encryption Solution

- **Solution:** We propose a Novel approach of TruStream[®] Computing – a **Non-Von Neumann Architecture** which includes:
 - A scalable fabric made up of programmable small cores on which an application developer can map dataflow parallelism and pipelining.
 - A standardized FHE API called the **Secure Computing Interface Framework (SCIFR API™)**
 - A series of control and data producer/consumer cores chained together to process the encrypted data together in coordinated fashion. This approach maps FHE algorithms dataflow into small cores chained together.

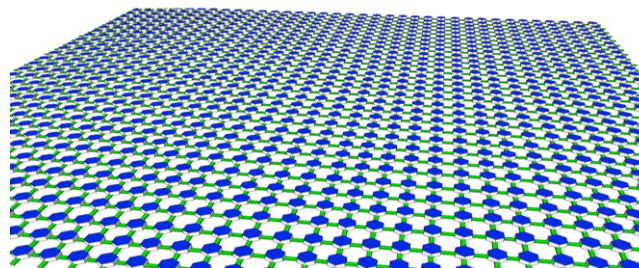
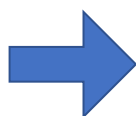
Scalable Factory for Data and Control



Source: <http://trickart.com/notes/wp-content/uploads/2014/09/faktur1140.gif>



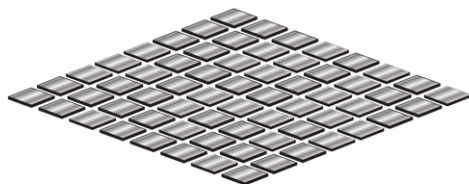
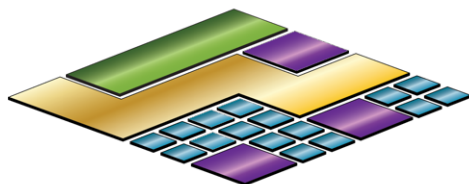
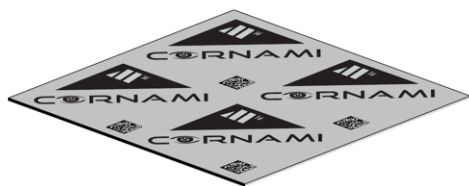
Today: Monolithic IC with Few Cores



Scalable number of cores that span many ICs

- ***Why are we still limited to a single monolithic chip for workload executions ?***
 - Load, Process and Store approach is performance hog for various reasons such as context switches, multi threads, non-deterministic software stack, cache dynamics, memory bottlenecks etc.
- ***Cornami solves all these problems***
 - Deterministic, seamlessly scalable and no context switch
 - Maximize concurrency and make full use of pipelining – Scalable performance, lower latency, lower silicon footprint and lower power use

Cornami Architecture Overview



TruStream® Reconfigurable Computing Fabric

Computational Fabric

Reconfigurable cores surrounded by SRAM and all integrated by a high-speed network. Fabric is linearly scalable across multiple chips, boards and systems to increase core for performance transparent to the software model.

Software Defined Cores

Created using FracTLcores® and under software control and configured based on function, for example:



FracTLcores®

Dynamically reconfigurable and independently programmable cores (each has its own scheduler) under software control to optimize performance within application programs.

TruStream®

- Software and programming environment supporting all 55 Berkeley patterns for parallelism
- Verified on numerous multi-core processors
- Large existing library of application software including Fully Homomorphic Encryption (FHE) and Privacy Preserving Machine Learning (PPML)

- To efficiently utilize and program a large number of resources in the Cornami fabric, we have developed the Secure Computing Interfacing Framework – SCIFR API™. Key features include
 - Coexist with all hardware resources such as CPU/GPU/FPGAs) including the scalable resources of a Cornami fabric.
 - Define multiple tiers of APIs for modules and operators
 - Polynomial Multiplications, Modulus, Bootstrapping, KeySwitching, NTTs, FFTs etc..
 - Enables seamless interfacing and integration with FHE Schemas and Frameworks



Secure Computing Interface Framework (SCIFR API™)

- **Allows complete execution of FHE circuits** on Cornami's scalable fabric.
 - Breaks free of the co-processor accelerator model by allowing an entire FHE circuit, FHE encryption, and FHE decryption to execute entirely on Cornami hardware
- Flexibility to achieve specific FHE application **throughput, latency, footprint, and/or performance goals** through compile time configuration
- Allows the adoption of new and emerging FHE algorithms
- Keeps FHE data movement within the Cornami scalable fabric allowing the stitching of multiple sub-topologies into a larger topology

SCIFR API™ – Abstraction Layers

Application	Algorithms and data
Independent	OpenAPIs Configuration Tool ML Conversion Tool ...
Framework	Seal OpenFHE HELib TFHE ...
Schema	FHE BFV CKKS TFHE ...
Hardware Interface	Calls to the Hardware Interface Layer (HIL) will find and leverage hardware assist Ability to load entire (logic/arithmetic) circuit for Cornami H/W execution
Operator	Modulus Key Switching Bootstrapping FFT NTT ...
Silicon	CPU GPU FPGA Cornami ...

- Cornami supports a scalable core fabric to enable FHE applications to execute at plaintext speeds.
 - Supports compile time configuration to achieve throughput, latency, footprint, and/or performance goals.
- Cornami's SCIFR API™ supports all hardware platforms and enables the execution of an entire FHE application circuit on a scalable computational fabric eliminating the need for a coprocessor approach
- Cornami's architecture is stream-based supporting a highly scalable and efficient core fabric that easily spans an arbitrary number of ICs.
- Cornami is committed to developing FHE standards and benchmarking to enable customers evaluate overall FHE circuit performance

