# Hierarchical Identify Verify Exploit (HIVE) Program

Trung Tran
DARPA/MTO

Proposers Day Brief
DARPA-BAA-16-52
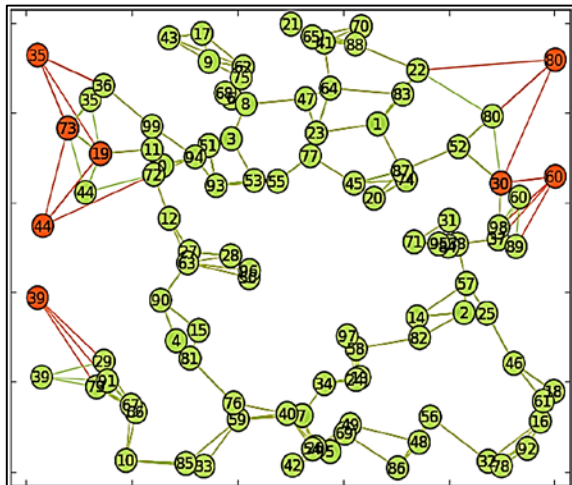
**HIVE will create a graph analytics processor that achieves 1000x improvement in processing efficiency**

- This will enable relationships between events to be discovered as they unfold in the field rather than relying on forensic analysis in data centers
- This will enable data scientists to make associations previously thought impractical due to the amount of processing required
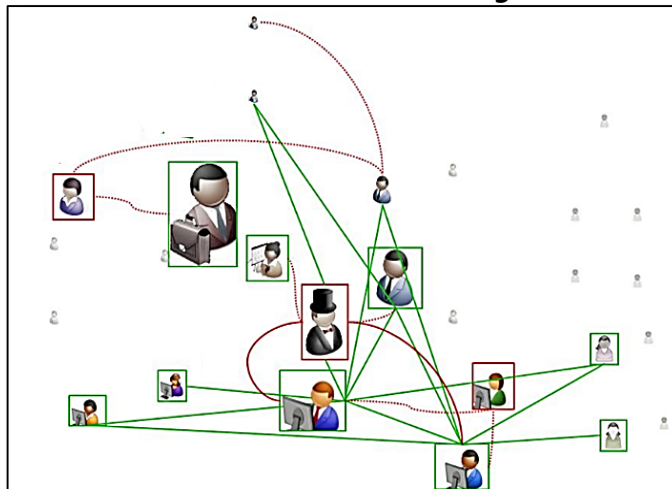
## Cyber Security



**Which cyber events are probes on the network?**

- Who are they probing and who have they infected in the network?

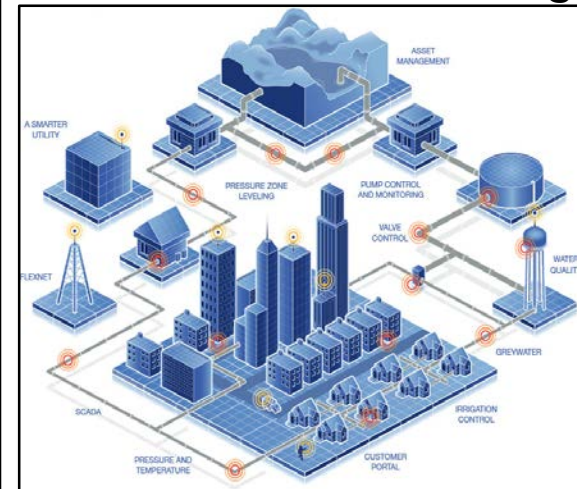- Only a small number of events are probes – graph is sparse.

## Social Media Analysis



**Who influences me to buy a product?**

- Who has access to my social media pages and what are they saying to me?

- Since only a few people have direct influence on me – graph is sparse.
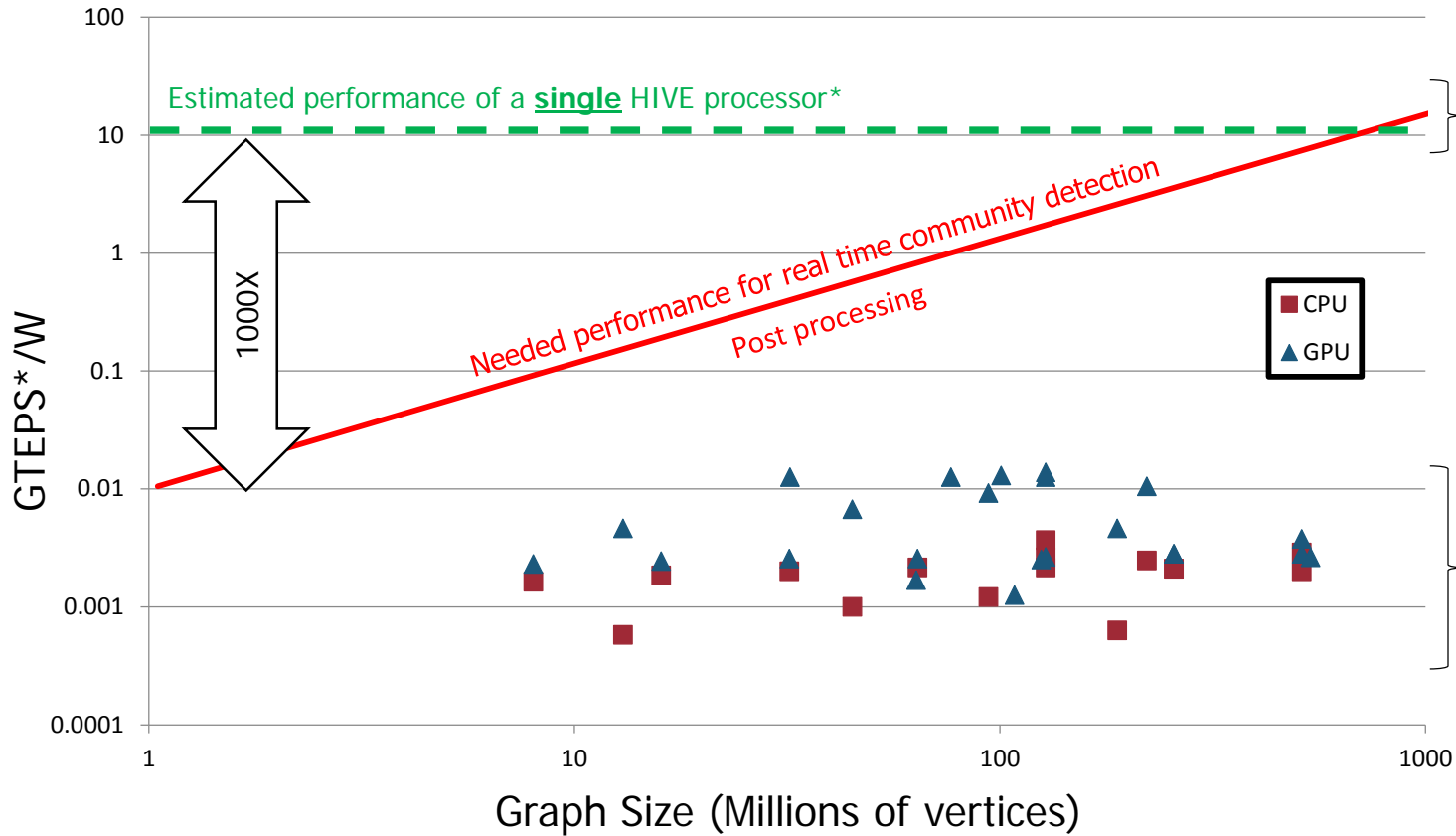
## Infrastructure Monitoring



**Can I spot failures before they become critical?**

- How do I avoid cascading failures and what are the system dependencies?

- Only a small number of critical dependencies – graph is sparse.

---

**Graph analytics is beginning to be applied to a broad set of problems**

# Graph analytics today requires large data centers



The HIVE program will develop a single processor capable of efficiently performing community detection on a graph of up to a billion vertices in real time.
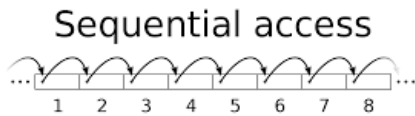
Current single chip CPU/GPU hardware cannot efficiently process large graphs in real time. To overcome processor limitations, large data centers are required.

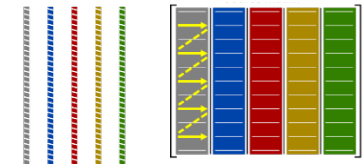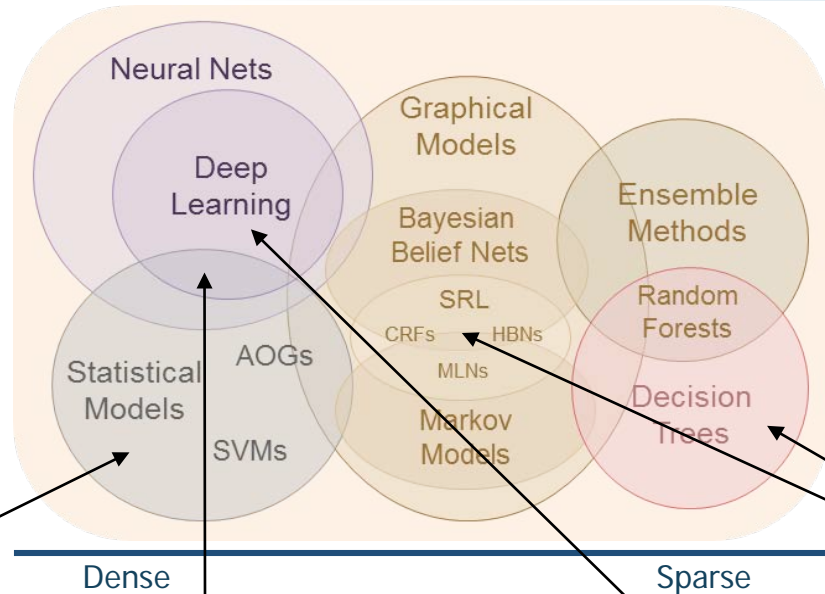HIVE aims to enable scalable, real-time graph analytics at the network edge

* GTEPS = Giga Traversed Edges Per Second

4

# HIVE – Today's hardware is focused on dense data

**Sequential access**

...1 2 3 4 5 6 7 8...

**Random access**

1 3 7 2 8 6 4 5

Sequential access is <u>good for dense data</u>
Sparse data requires random access



Neural Nets

Deep Learning

Graphical Models

Bayesian Belief Nets

Ensemble Methods

SRL

CRFs     HBNs

MLNs

Statistical Models

AOGs

SVMs

Random Forests

Decision Trees

Markov Models

Dense → Sparse

<u>Lower level primitives</u>
(5x5 Matrix)
- 25 Scalar operations
- 5 Vector operations
- 1 Matrix operation

## Intel CPU
- Sequential processing
- Sequential memory access
- Slow (20GB/s) to memory
- Limited scalability (16GB/s)
- Optimized for Statistics

Source: Intel

## Nvidia GPU
- Parallel processing
- Sequential memory access
- Faster (288GB/s) to memory
- Limited scalability (20GB/s)
- Used for CNNs

Source: Nvidia

## Google TPU
- Parallel processing
- Sequential memory access
- Slow (20GB/s) to memory
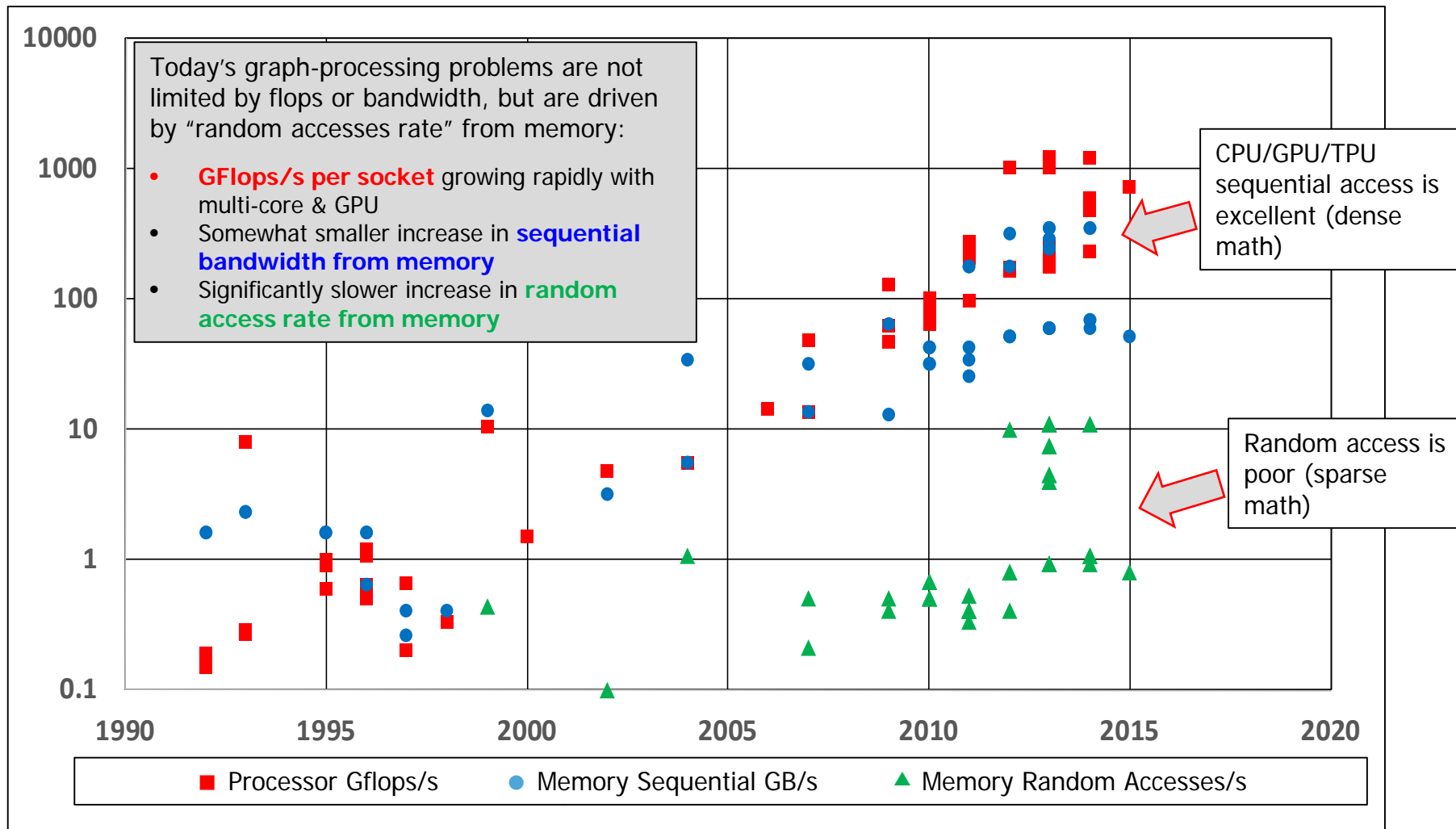- Limited scalability (16GB/s)
- Optimized for DNNs

Source: Google

## HIVE
- Parallel processing
- Parallel memory access
- Fastest (TB/s) to memory
- Higher scalability (TB/s)
- Optimized for Graphs

Today's graph-processing problems are not limited by flops or bandwidth, but are driven by "random accesses rate" from memory:

- **GFlops/s per socket** growing rapidly with multi-core & GPU
- Somewhat smaller increase in **sequential bandwidth from memory**
- Significantly slower increase in **random access rate from memory**

CPU/GPU/TPU sequential access is excellent (dense math)

Random access is poor (sparse math)

Legend:
- ■ Processor Gflops/s
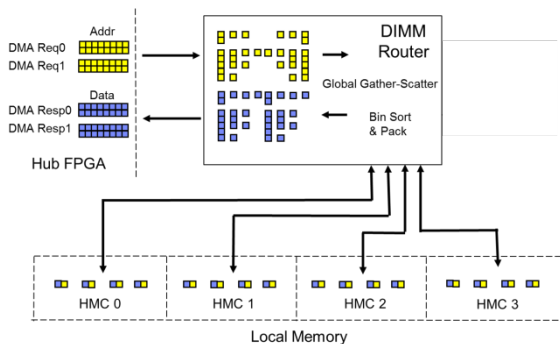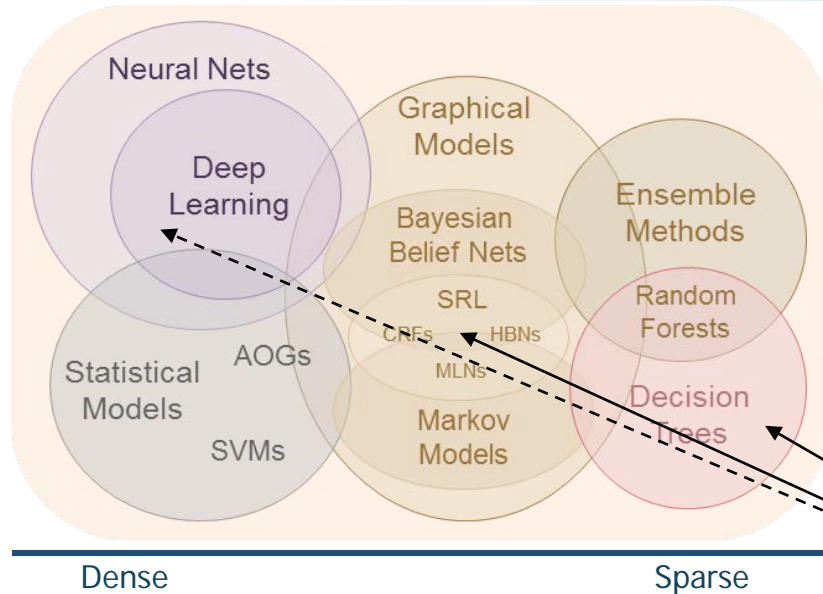- ● Memory Sequential GB/s
- ▲ Memory Random Accesses/s
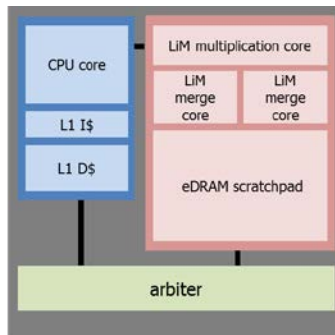
\* EMU Technologies Design Review.  March 2016

# HIVE is focused on sparse data

**New Approach**

Create a new hardware architecture optimized for sparse graphical models and decision trees



Dense ➝ Sparse



Parallel memory access enables random access



**Graph Primitives**
SpRef/SpAsgn, SpGEMM, SPEWiseX, Scale, Reduce, and Apply

Matrix data format/operations enables more efficient processing



**HIVE**
- Parallel processing
- Parallel memory access
- Fastest (TB/s) to memory
- Higher scalability (TB/s)
- Optimized for Graphs

Large graphs → Represented as matrices → Random access | SpRef/SpAsgn, SpGEMM, SPEWiseX, Scale, Reduce, and Apply | Graph primitives → HIVE

| Data mapping | Graph analytics tool packages | New memory architectures | Graph primitives | System on chip design flows |

SpRef/SpAsgn, SpGEMM, SPEWiseX, Scale, Reduce, and Apply

| | Define graph primitives | Create data format model | Define data flow model |
|---|---|---|---|
| | ⬇ | ⬇ | ⬇ |
| **Graph Software**<br><br>What should be accelerated? | Define linear algebra building blocks which can be accelerated | Map graph matrix into subarrays which can allow for easy memory mapping | Define data movement from processor to memory and between processors |
| | Accelerators | Memory | Scaling |
| **Graph Accelerator**<br><br>How should it be accelerated? | Develop hardware accelerators for each building block | Create memory controller which optimizes data movement based on mapping | Develop bus architectures to avoid congestion in data movement |

# HIVE – Program structure



Large graphs → Represented as matrices → Random access → SpRef/SpAsgn, SpGEMM, SPEWiseX, Scale, Reduce, and Apply — Graph primitives → HIVE

- Targeting graph applications
- Dense and sparse data
- Built for scaling
- Intended for embedded

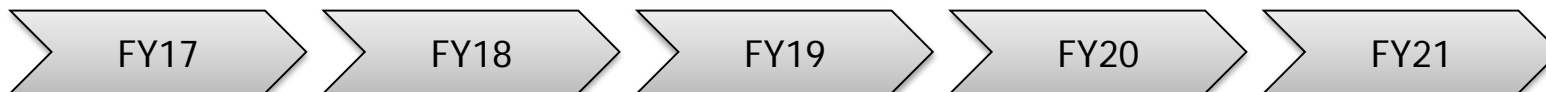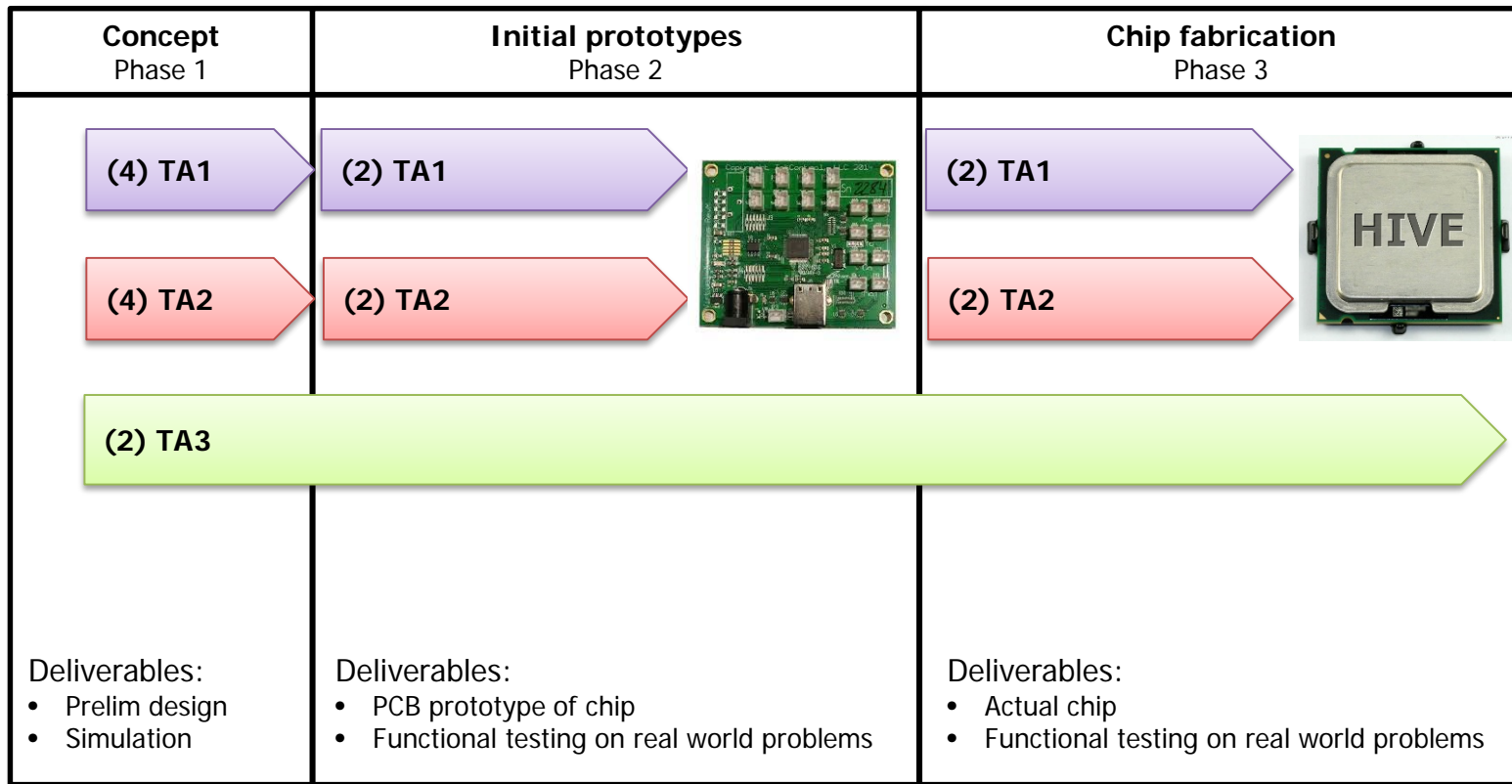| Challenge problem areas | TA 2: Graph analytics toolkits (ref: Tensorflow, CUDA) | TA 1: Graph analytics accelerator (ref: TPU, GPU) | Evaluation framework |
|---|---|---|---|
| | • Enable real-time streaming graph analytics<br>• Designed for hardware acceleration<br>• Generally applicable across a series of graph problems | • Runs at <20W<br>• Enable reduced data movement/processor idle time to <50%<br>• Allow for 95% memory BW efficiency at 100% random access (R/W) | 100X improvement |

**TA3:** Evaluator

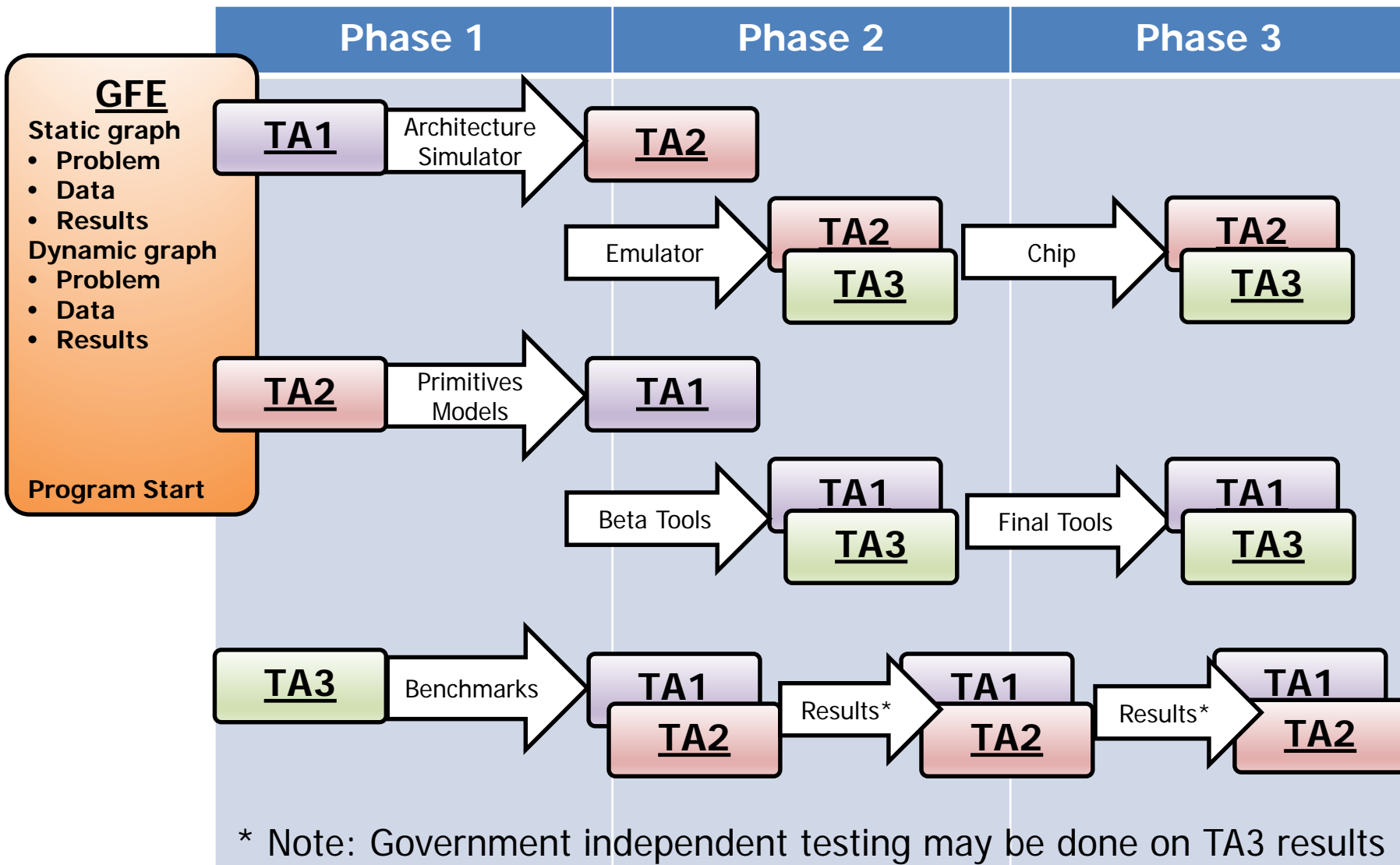Blog Analysis — Community thought leaders

FaceBook - 300 M users — Community Activities

Semantic Web — People, Places, & Actions

SmartGrid — N-x contingency analysis

| Concept<br>Phase 1 | Initial prototypes<br>Phase 2 | Chip fabrication<br>Phase 3 |
|---|---|---|
| (4) TA1 | (2) TA1 | (2) TA1 |
| (4) TA2 | (2) TA2 | (2) TA2 |
| (2) TA3 | | |
| Deliverables:<br>• Prelim design<br>• Simulation | Deliverables:<br>• PCB prototype of chip<br>• Functional testing on real world problems | Deliverables:<br>• Actual chip<br>• Functional testing on real world problems |

FY17 &gt; FY18 &gt; FY19 &gt; FY20 &gt; FY21

www.darpa.mil