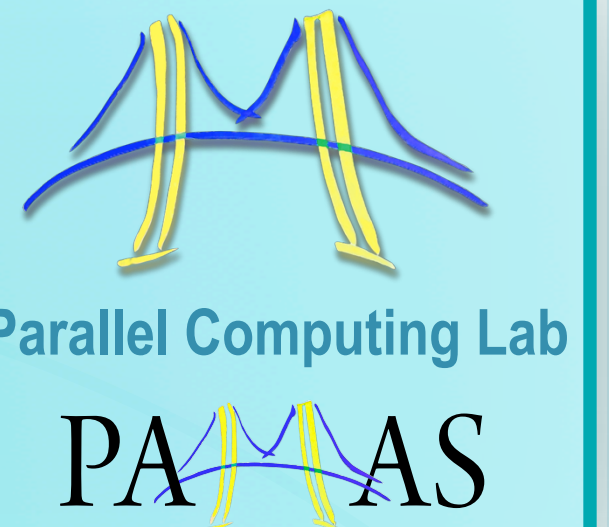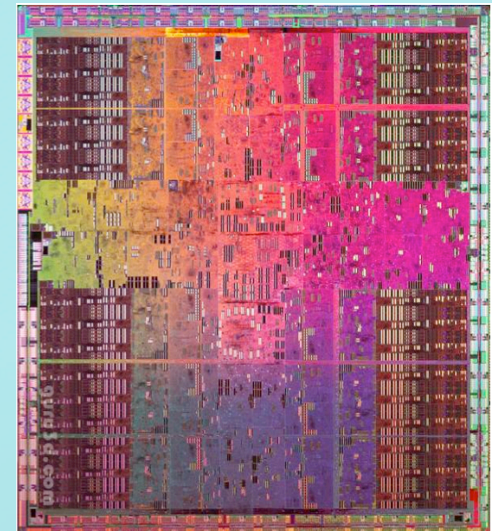# Application-level Trade-offs for WFST-based Large Vocabulary Continuous Speech Recognition on a Graphics Processing Unit

Jike Chong, Ekaterina Gonina, Youngmin Yi, Kurt Keutzer, Department of Electrical Engineering and Computer Science, University of California, Berkeley

Parallel Computing Lab
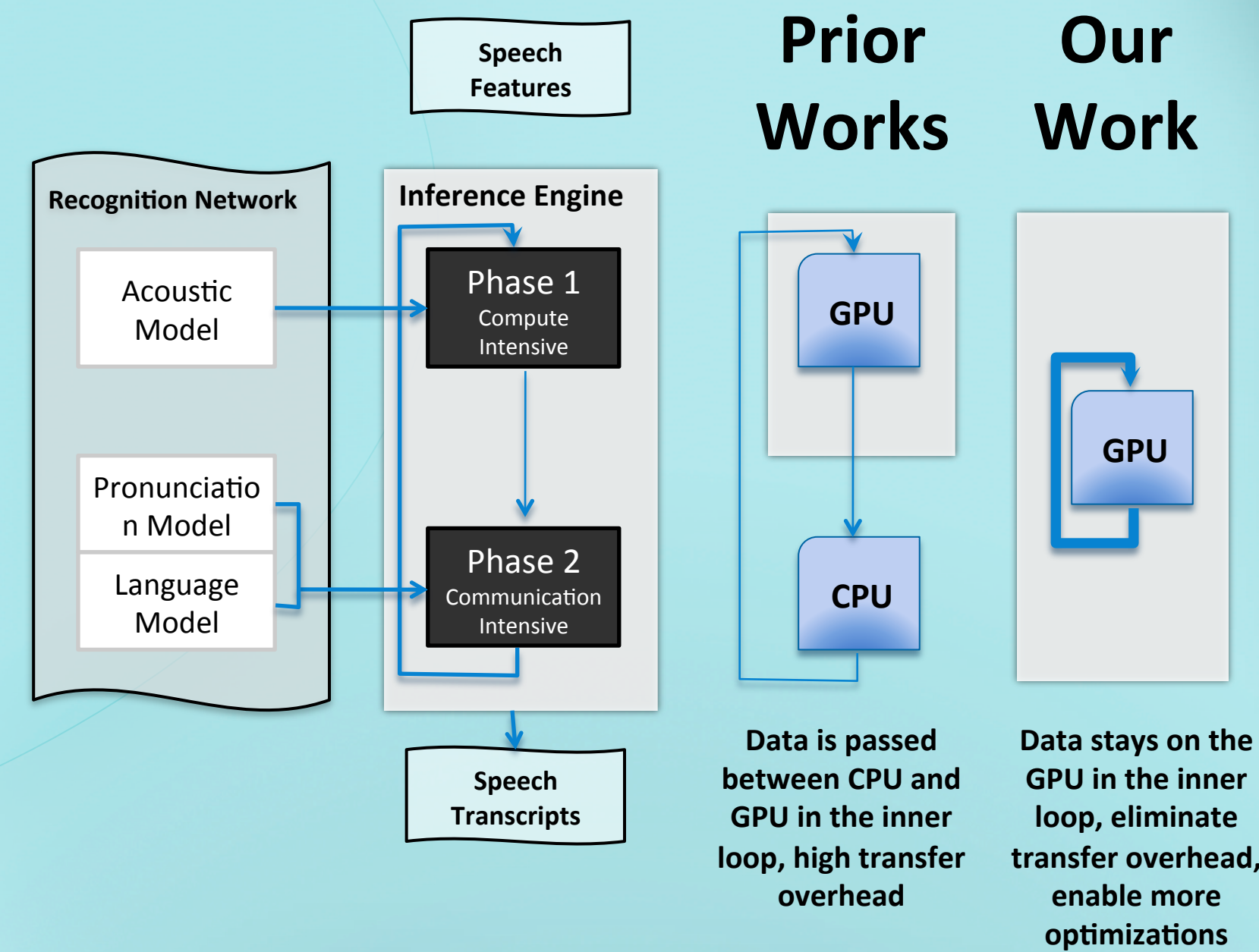PALLAS

## Emerging Manycore Platforms

**NVIDIA GTX280 30 cores**

- Architecture trend:
  - Increasing vector unit width
  - Increasing numbers of cores per die
- Application implications:
  - Must optimize synchronization cost
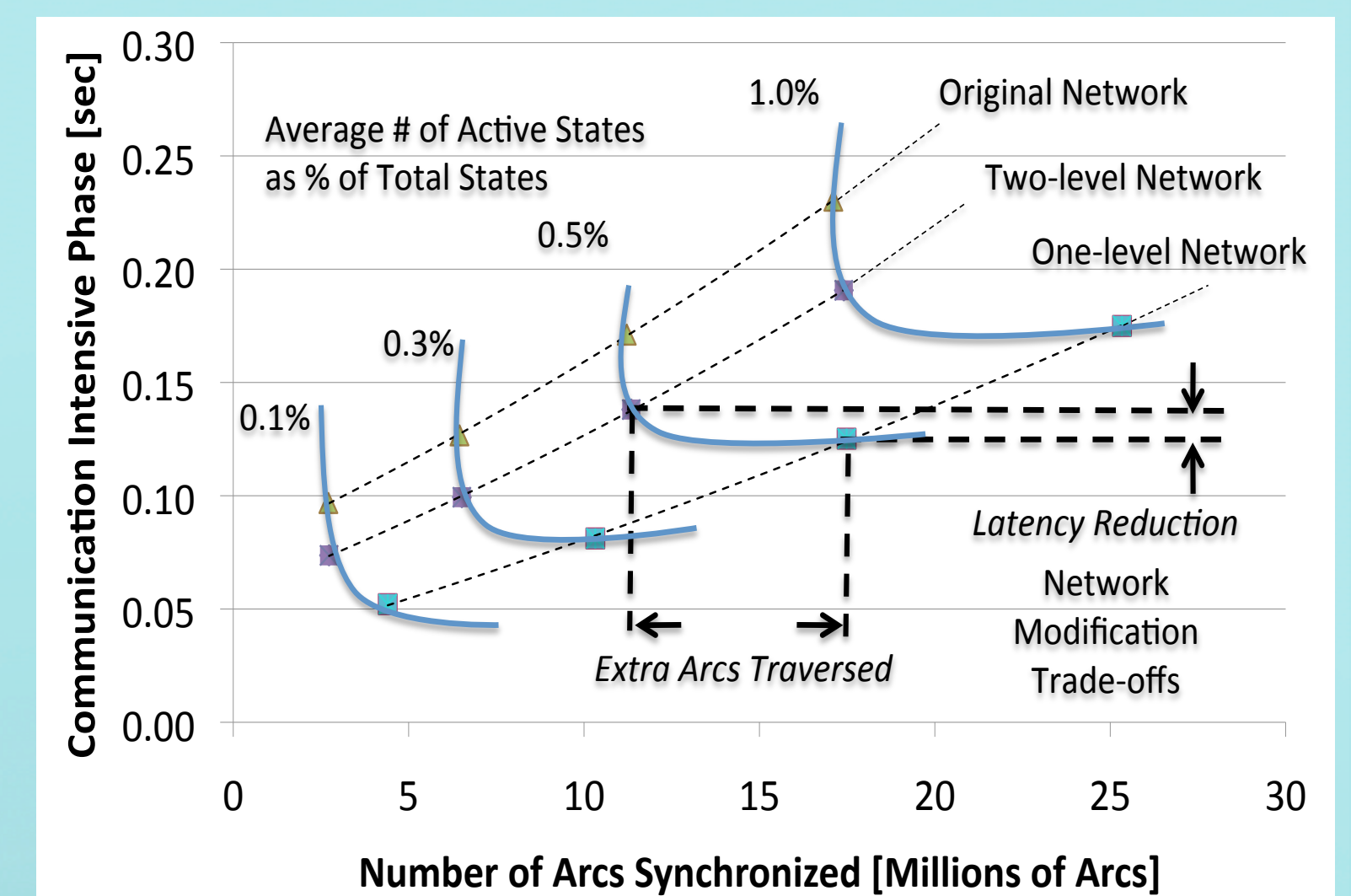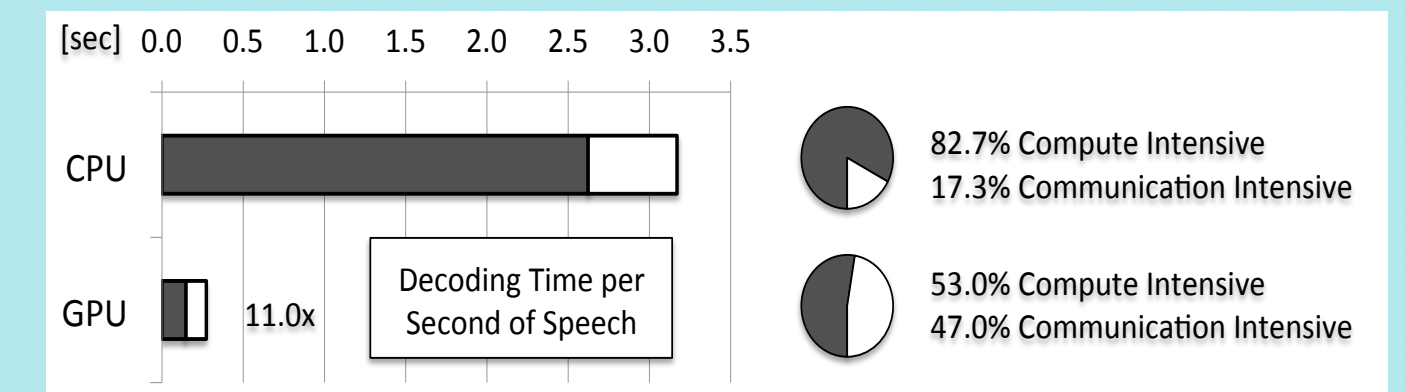  - Must increase SIMD efficiency

Ongoing work investigates algorithm design space that optimizes for data parallel manycore programming.

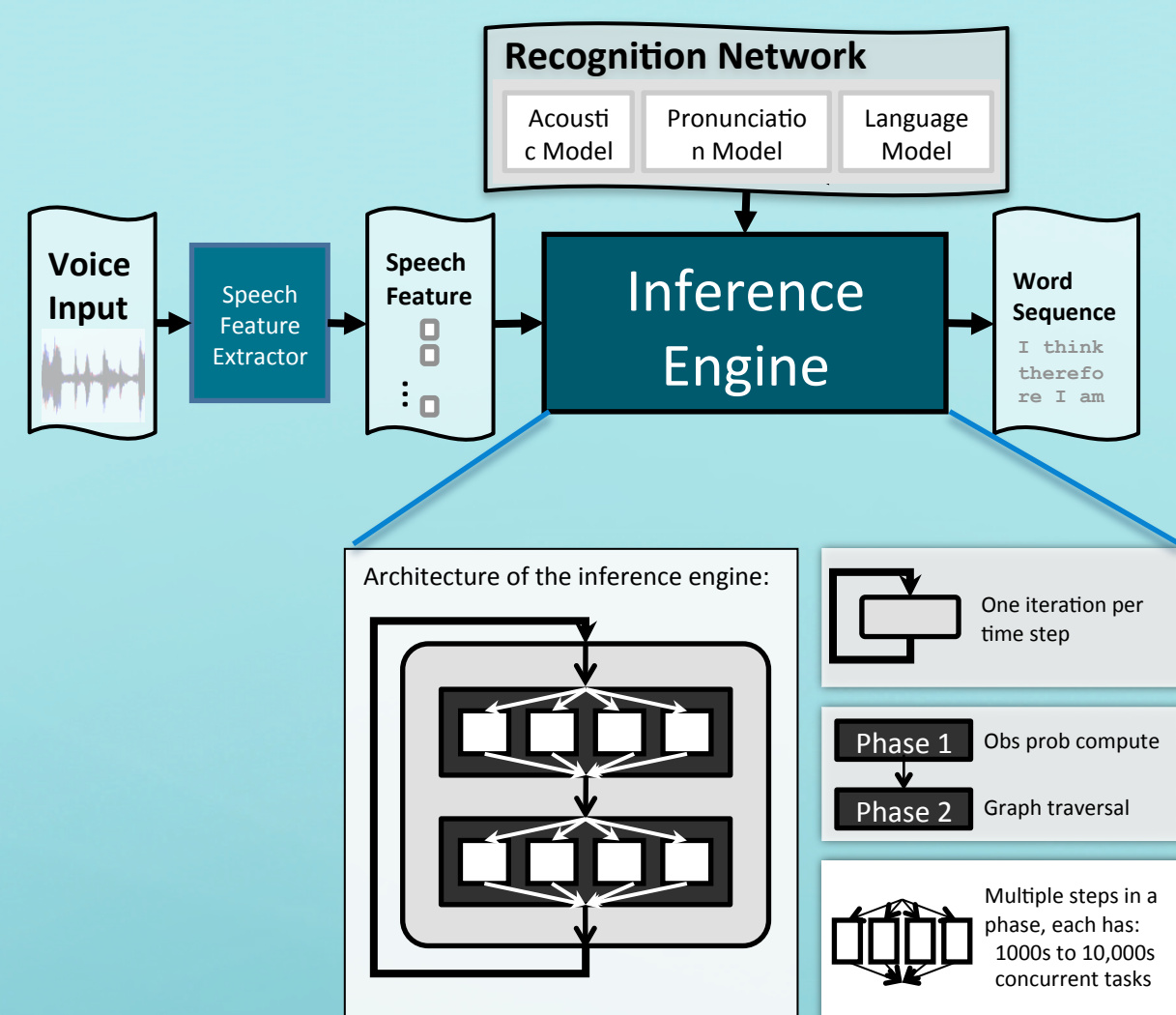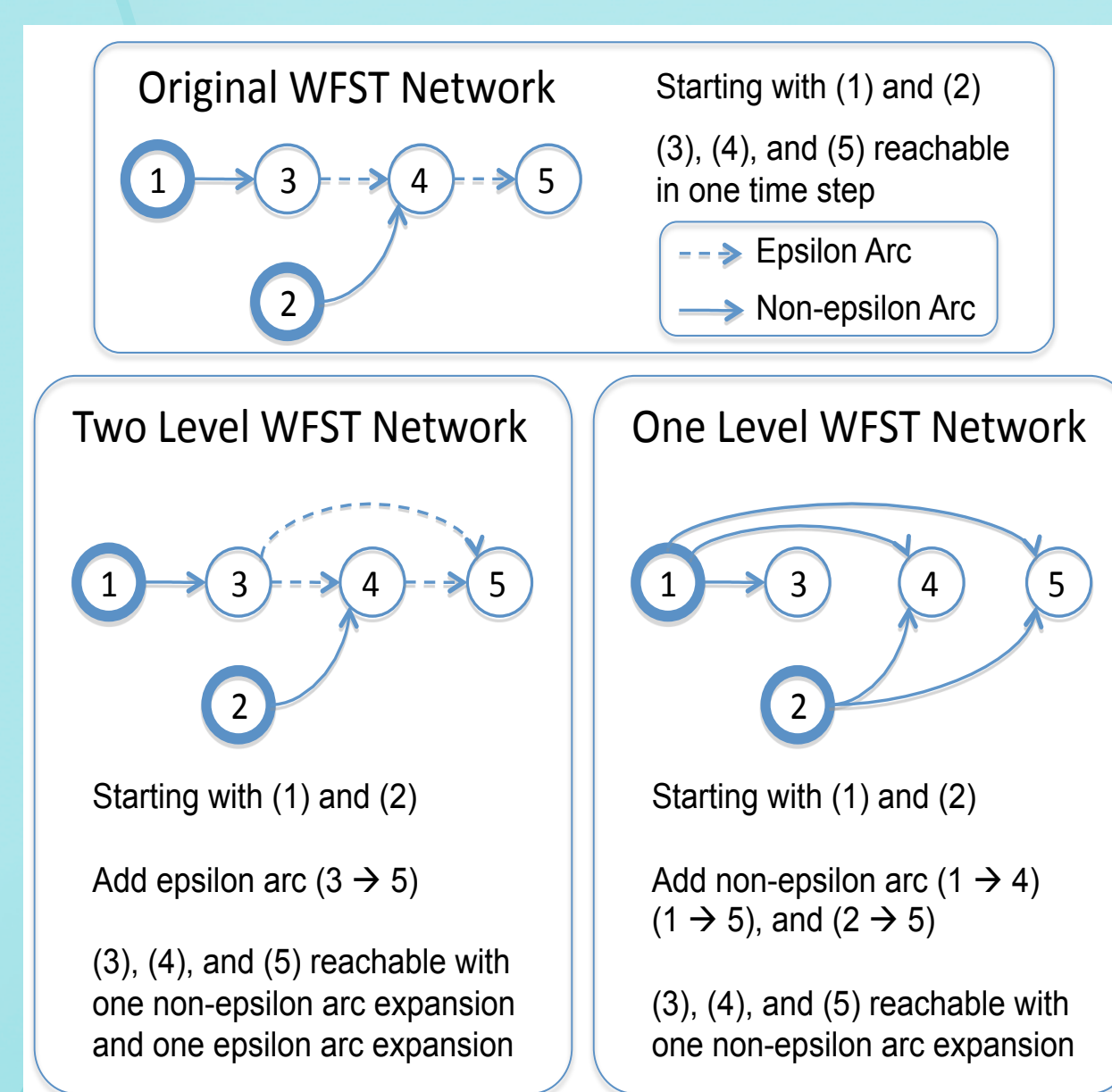## Speech Recognition Inference Engine Characteristics

- Parallel graph traversal through irregular network
  - Guided by a sequence of input audio vectors
  - Computing on continuously changing data working set
- Implementation challenges
  - Define a scalable software architecture to expose fine-grained application concurrency
  - Efficiently synchronize between an increasing number of concurrent tasks
  - Effectively utilize the SIMD-level parallelism

**Recognition Network**
Acoustic Model | Pronunciation Model | Language Model

Voice Input → Speech Feature Extractor → Speech Feature → Inference Engine → Word Sequence

I think therefore I am

Architecture of the inference engine:
- One iteration per time step
- Phase 1 — Obs prob compute
- Phase 2 — Graph traversal
- Multiple steps in a phase, each has: 1000s to 10,000s concurrent tasks

## System Architecture

Speech Features

**Recognition Network**
Acoustic Model
Pronunciation Model
Language Model

**Inference Engine**
Phase 1 — Compute Intensive
Phase 2 — Communication Intensive

Speech Transcripts

**Prior Works**
GPU
CPU
Data is passed between CPU and GPU in the inner loop, high transfer overhead

**Our Work**
GPU
Data stays on the GPU in the inner loop, eliminate transfer overhead, enable more optimizations

## Algorithm Design Space Exploration

**Original WFST Network**
Starting with (1) and (2)
(3), (4), and (5) reachable in one time step
- - - Epsilon Arc
→ Non-epsilon Arc

**Two Level WFST Network**
Starting with (1) and (2)
Add epsilon arc (3 → 5)
(3), (4), and (5) reachable with one non-epsilon arc expansion and one epsilon arc expansion

**One Level WFST Network**
Starting with (1) and (2)
Add non-epsilon arc (1 → 4) (1 → 5), and (2 → 5)
(3), (4), and (5) reachable with one non-epsilon arc expansion

- Explore efficient graph traversal technique
  - Vary the amount of flattening of the WFST network to two levels or one level
  - The flattening increases the number of arcs to traverse in the algorithm

## Evaluation of the Inference Engine

[sec] 0.0 0.5 1.0 1.5 2.0 2.5 3.0 3.5
CPU
GPU 11.0x — Decoding Time per Second of Speech

82.7% Compute Intensive
17.3% Communication Intensive

53.0% Compute Intensive
47.0% Communication Intensive

- Speed up varies btw phases:
  - 4-18x for compute intensive phases
  - 3-4x for communication intensive phases
  - Communication intensive phases becoming proportionally more important

- Speedup for phase 2:
  - Moving to 2-level WFST network provides 17-24% speed up with minimal increase in arcs traversed
  - Moving to one level WFST network provides an additional 8-29% speedup at the expense of traversing 48-62% more arcs
  - Less than 8% sequential overhead

Communication Intensive Phase [sec] (y-axis: 0.00 to 0.30)
Number of Arcs Synchronized [Millions of Arcs] (x-axis: 0 to 30)

Average # of Active States as % of Total States
1.0%, 0.5%, 0.3%, 0.1%

Original Network
Two-level Network
One-level Network

Network Modification Trade-offs
Latency Reduction
Extra Arcs Traversed

## Conclusions

- Defined and implemented a parallel software architecture:
  - 5-8% sequential overhead
  - Significant potential for further speedup in future platforms
- Implemented the entire inference engine on the GPU
  - Both GMM computation and graph traversal phase are implemented in data-parallel routines
- Explored the algorithmic design space for WFST network optimization for data parallel operations
  - Network flattening critical for efficient data parallel operation

We expect that an efficient speech recognition engine will be a key component in many exciting new applications to come!