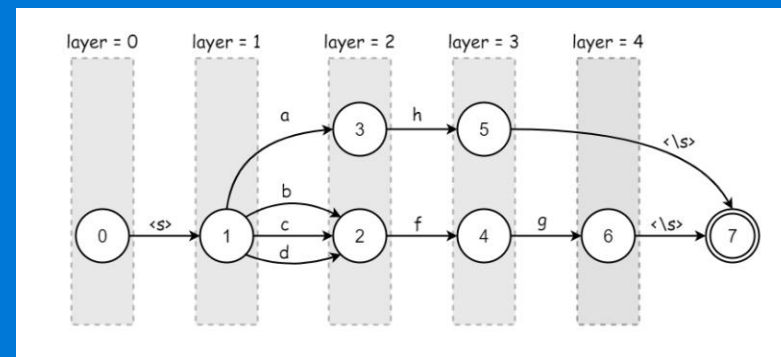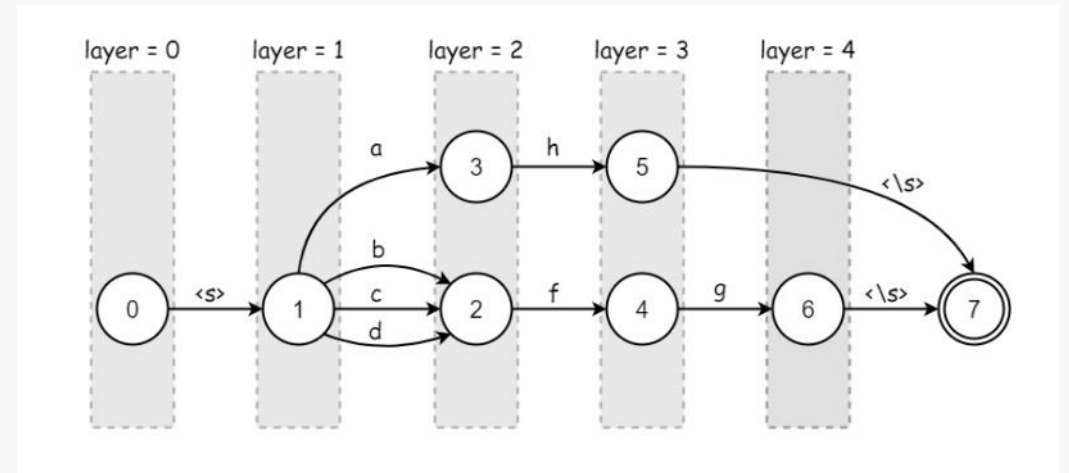# A GPU-based WFST Decoder with Exact Lattice Generation

**Zhehuai Chen**, Justin Luitjens, **Hainan Xu, Yiming Wang**,

Daniel Povey, Sanjeev Khudanpur

The Center For Language and Speech Processing at the Johns Hopkins University

SJTU SPEECH LAB
上海交通大学智能语音实验室

NVIDIA.
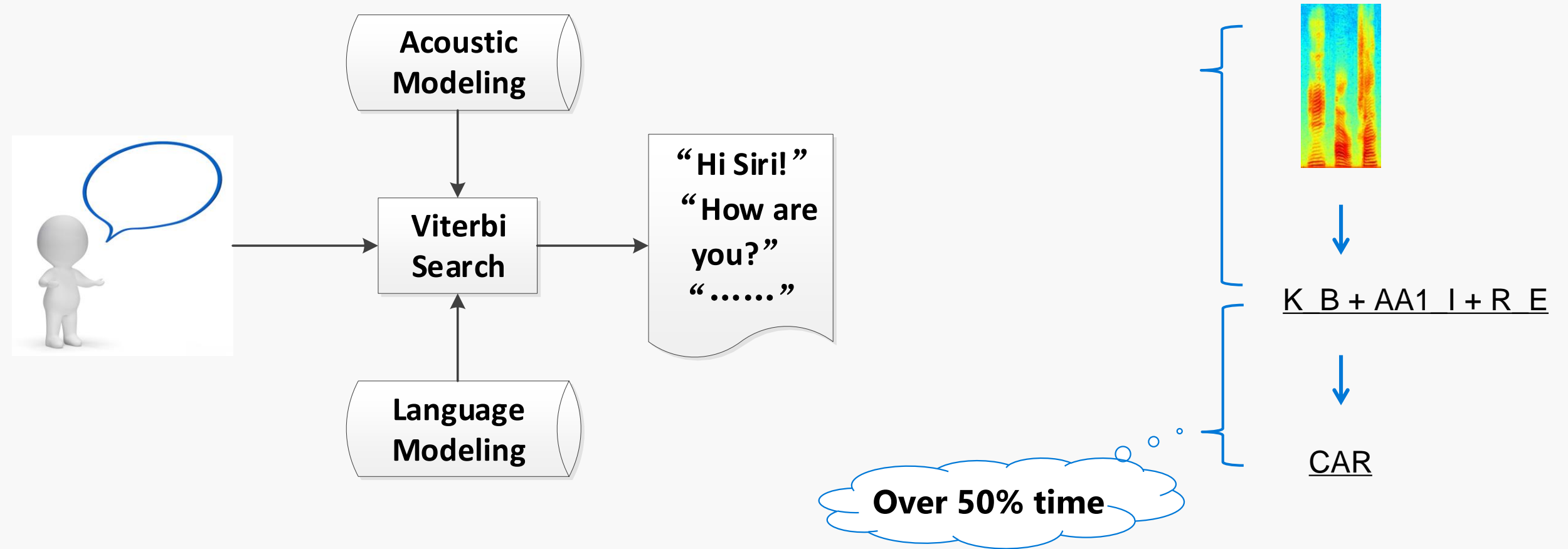
KALDI

# Outline

- Introduction
  - ASR and WFST Decoding
  - GPU-based Parallel Computing

- <u>Parallel Viterbi Decoding[1]</u>
  - Framework
  - Token Recombination
  - Load Balancing
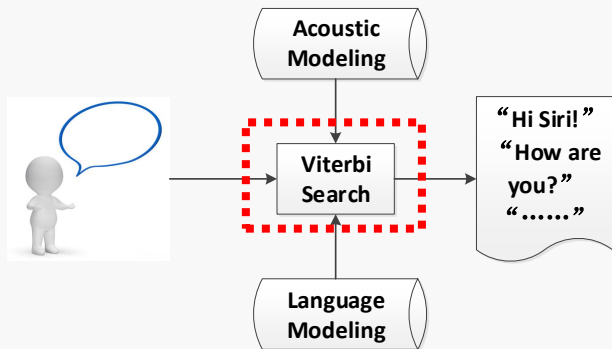  - Lattice Processing

- Experiments



1. https://github.com/chenzhehuai/kaldi/tree/gpu-decoder

# Introduction

- Automatic Speech Recognition (ASR)

# Introduction

- Weighted finite state transducer (WFST) Decoding



| word | Pronunciation | Score[1] |
|------|---------------|-------|
| <s> | sil | 0.1 |
| </s> | sil | 0.6 |
| START | s t aa r t | 0.5 |
| STOP | s t aa p | 0.4 |
| IT | ih t | 0.3 |

1. There is a score for each arc. For simplicity, we do not draw them.

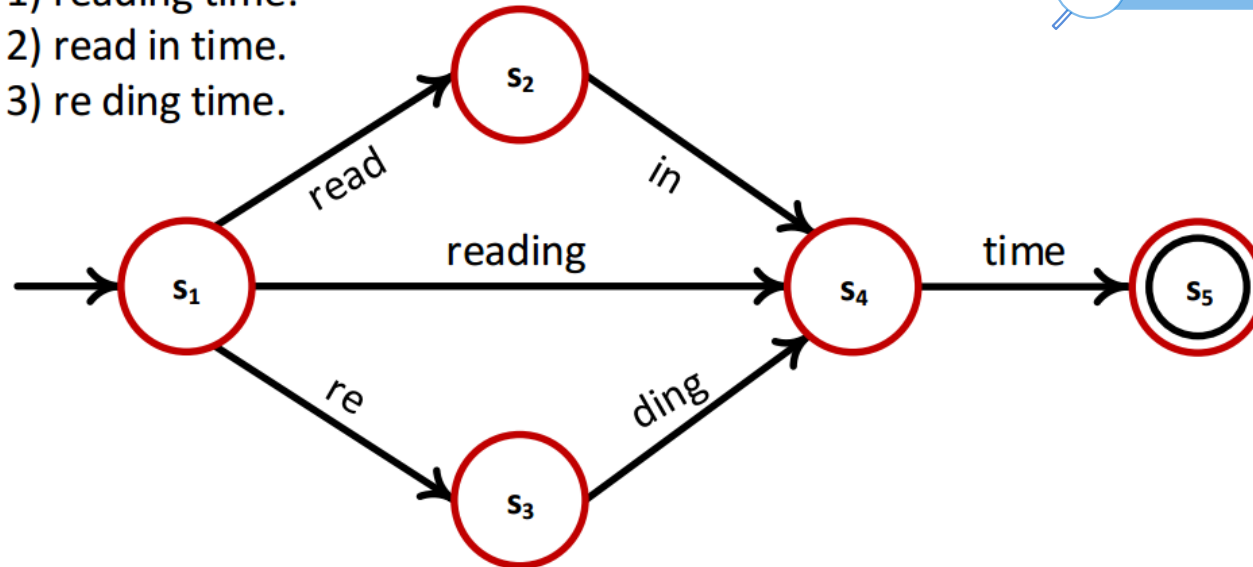# Introduction

- Lattice Processing

# Introduction

- ## GPU-based Parallel Computing
  - ### Matrix calculation
  - ### Sequence batching



**Training**



**Inference**

# Introduction

- This Work: GPU-based Decoding



Parallel Viterbi Search

# Introduction

- This Work: GPU-based Decoding



Kaldi ASR pipeline

audio → feature extraction → features → acoustic modeling → posteriors → WFST decoding (CPU) → text/lattice

Kaldi ASR pipeline
(improved)

audio → feature extraction → features → acoustic modeling → posteriors → WFST decoding (GPU) → text/lattice

# Introduction

- This Work: GPU-based Decoding



- GPU memory is enough for 1-pass WFST

# Framework

- Parallel Viterbi decoding
  - **Future:** out-going arcs, e.g. b, c & d

# Framework

- ## Parallel Viterbi decoding
  - **Future:** out-going arcs, e.g. b, c & d
  - **History:** per layer, e.g. 3 & 2

# Framework

- ## Parallel Viterbi decoding
  - **Future:** out-going arcs, e.g. b, c & d
  - **History:** per layer, e.g. 3 & 2
  - **Utterance:** decoding with separate GPU kernels

# Framework

- ## System overview
  - ### 2 GPU streams & 1 CPU thread

# 1ˢᵗ Problem: Token recombination

- Token recombination

**Language Model**

**Dictionary**

**DNN-HMM**

$$\max \; P(\mathbf{W}|\mathbf{O}) \quad \propto \quad \max \; P(\mathbf{W}) \cdot P(\mathbf{L}|\mathbf{W}) \cdot \sum_{\mathbf{q} \in \mathcal{A}(\mathbf{L})} p(\mathbf{O}, \mathbf{q}|\mathbf{L}) \qquad (1)$$

$$\propto \quad \max \; P(\mathbf{W}) \cdot P(\mathbf{L}|\mathbf{W}) \cdot \max_{\mathbf{q} \in \mathcal{A}(\mathbf{L})} p(\mathbf{O}, \mathbf{q}|\mathbf{L}) \qquad (2)$$

# Token recombination

- Token recombination

$$\max P(\mathbf{W}|\mathbf{O}) \propto \boxed{\max} P(\mathbf{W}) \cdot P(\mathbf{L}|\mathbf{W}) \cdot \boxed{\max_{\mathbf{q} \in \mathcal{A}(\mathbf{L})}} p(\mathbf{O}, \mathbf{q}|\mathbf{L}) \qquad (2)$$

**Serial** Viterbi Search

We need a
**Thread Sync.**

**Parallel** Viterbi Search

# Token recombination



- ## Naïve implementation:
  - Critical section

```
1  while(*cur_tokv < next_tok) {     // check if we need to recombine token
2      acquire_semaphore (int*)&params.token_locks[nextstate]); // per token lock
3      if(*cur_tokv < next_tok) { // double check if we are min
4          *cur_tok=next_tok; // this is what we really want to do
5      }
6      release_semaphore (int*)&params.token_locks[nextstate]); // release per token lock
7      break; // exit loop as our update is done
8  } // end while
```

- Thread deadlock in earlier GPU architectures
- Slow because of **while loop** and **semaphore acquiring**

# Token recombination

- ## Proposed implementation:
  - Single atomic operation[1]



```
int64* atomicMax(*address, val) {
    original = *address;
    *address = max(val, *address);
    return original;
}
```

1. *atomicMax(*address, val)*. Computes *max(*address, val)*, writes the result to *address*, and returns the original *address* .

# Token recombination

- ## Proposed implementation:
  - Single atomic operation[1]



```
int64* atomicMax(*address, val) {
    original = *address;
    *address = max(val, *address);
    return original;
}
```

a

score=0.7
arc_id=1

| 0.7 | 1 |

1

b

score=0.8
arc_id=2

2

e

c

score=0.6
arc_id=3

3

| 0 | -1 |

4

score=0.9
arc_id=4

d

1. *atomicMax(*address, val)*. Computes *max(*address, val)*, writes the result to *address*, and returns the original *address* .

# Token recombination

- Proposed implementation:
  - Single atomic operation[1]

```
int64* atomicMax(*address, val) {
    original = *address;
    *address = max(val, *address);
    return original;
}
```



1. *atomicMax(\*address, val)*. Computes *max(\*address, val)*, writes the result to *address*, and returns the original *\*address* .

# Token recombination

- ## Proposed implementation:
  - ### Single atomic operation[1]

```
int64* atomicMax(*address, val) {
    original = *address;
    *address = max(val, *address);
    return original;
}
```

a

score=0.7
arc_id=1

1

b

score=0.8
arc_id=2

2

| 0.6 | 3 |
|-----|---|

score=0.6
arc_id=3

3

c

e

| 0.8 | 2 |
|-----|---|

4

score=0.9
arc_id=4

d

1. *atomicMax(*address, val)*. Computes *max(*address, val)*, writes the result to *address*, and returns the original *address* .

# Token recombination

- ## Proposed implementation:
  - Single atomic operation[1]



```
int64* atomicMax(*address, val) {
    original = *address;
    *address = max(val, *address);
    return original;
}
```
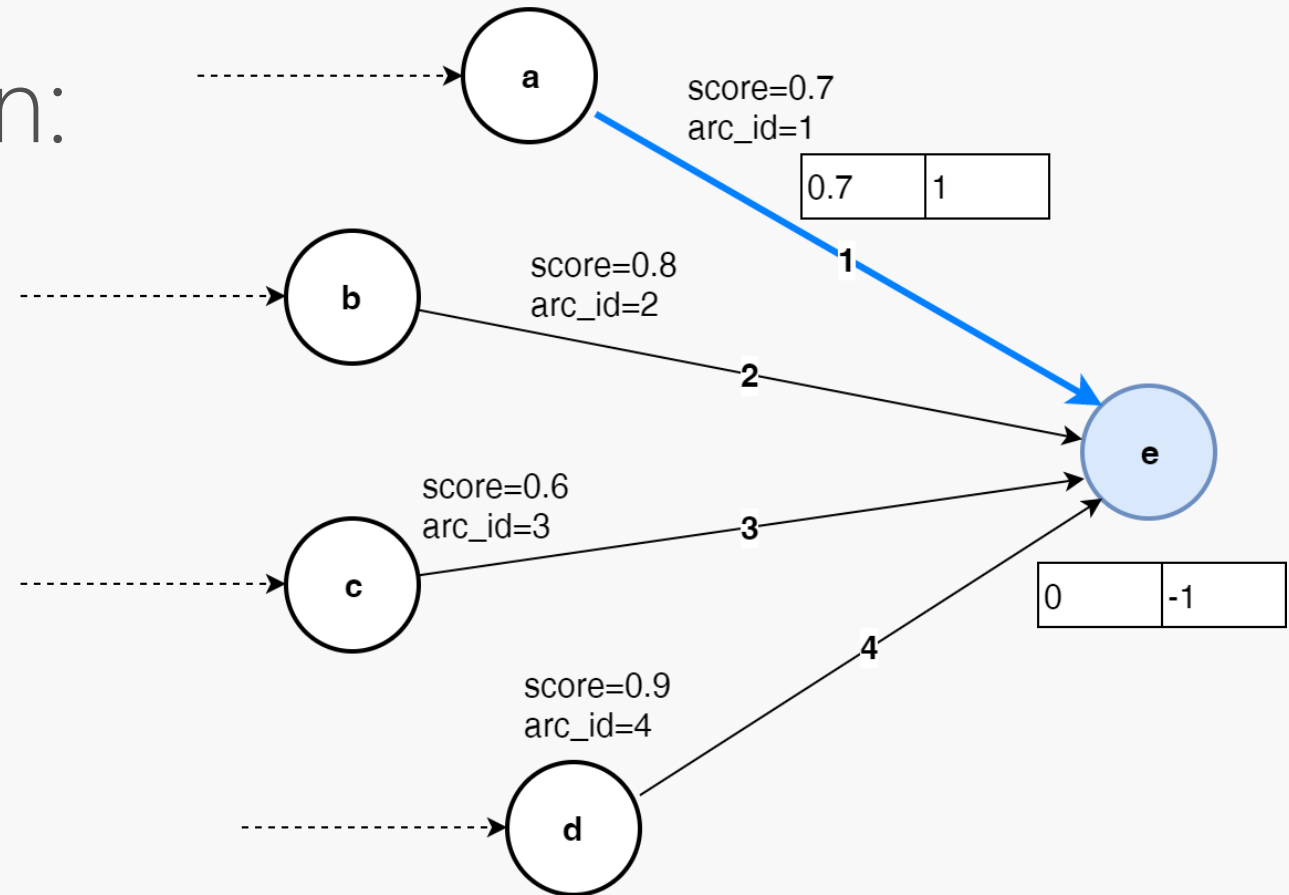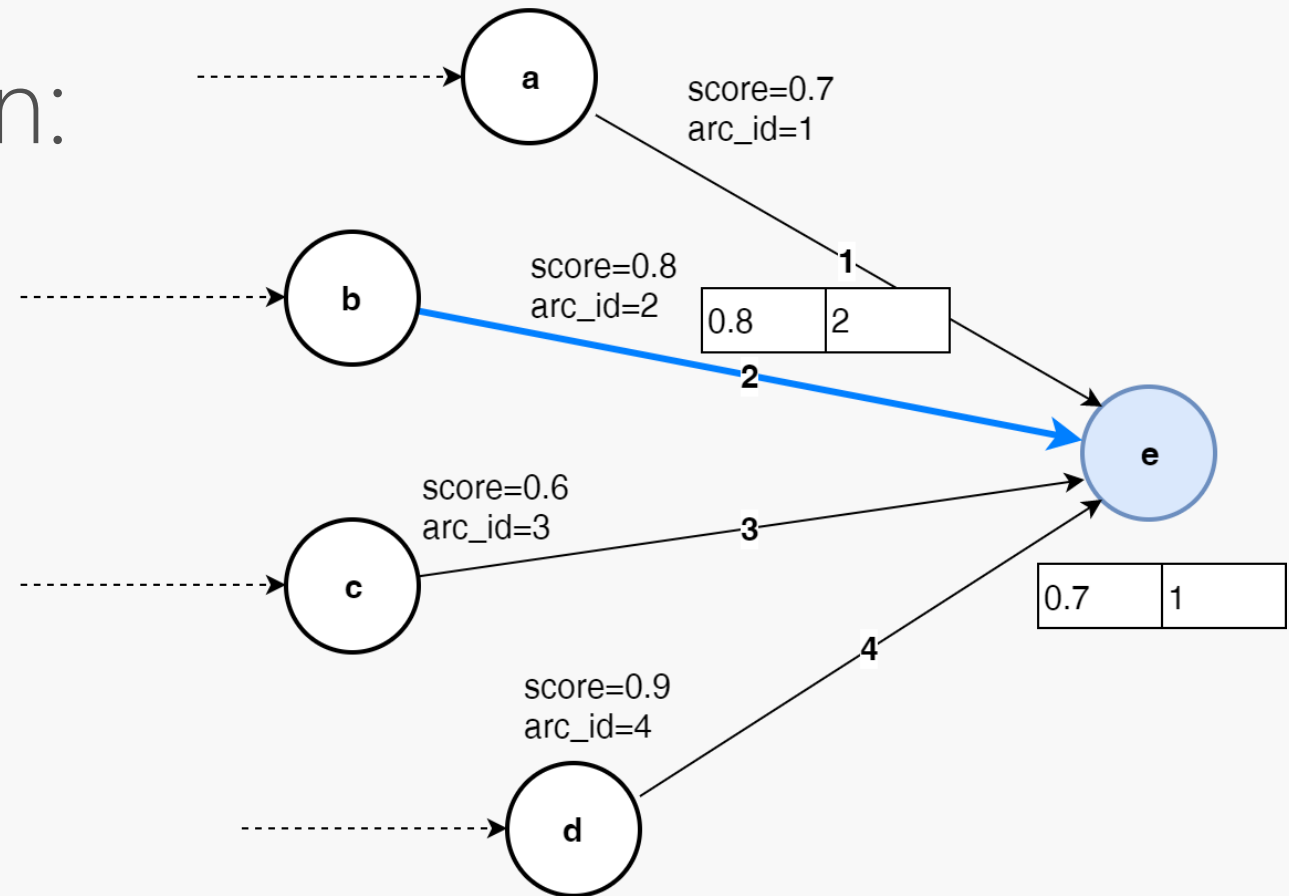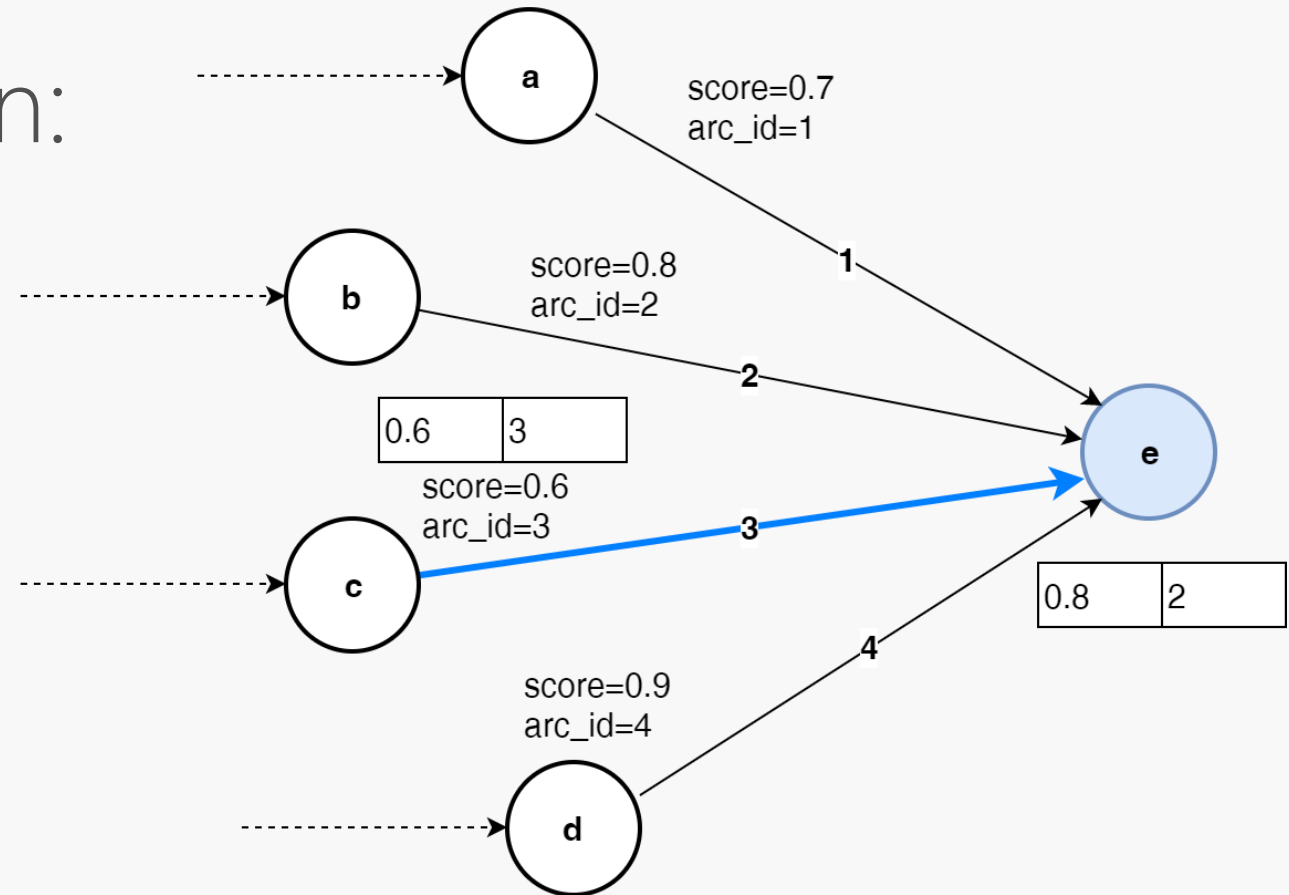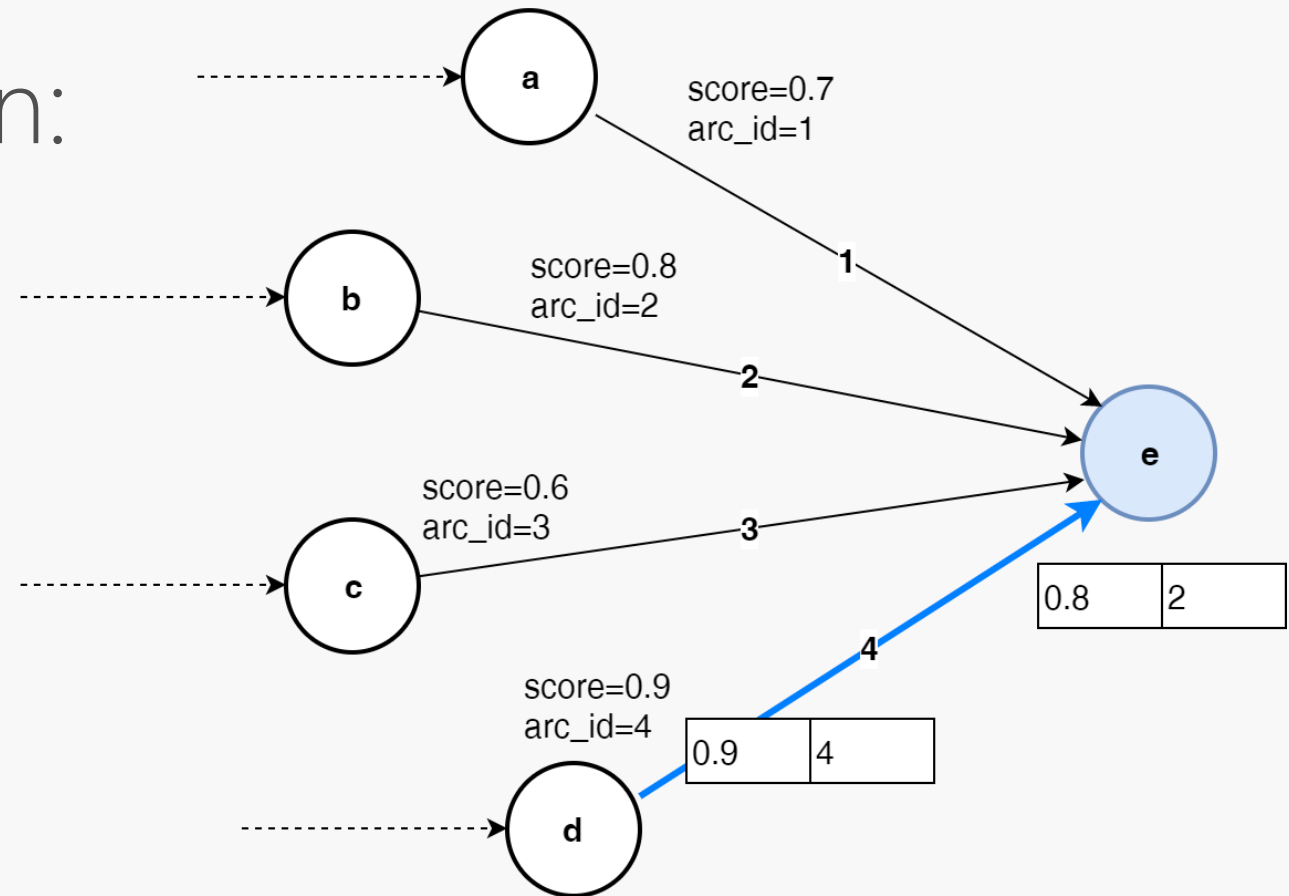
1. *atomicMax*(**address*, *val*). Computes *max*(**address*, *val*), writes the result to *address*, and returns the original **address* .

# Token recombination

- Proposed implementation:
  - Single atomic operation[1]



**Algorithm 1** Thread-level Token Recombination (Inputs: accumulated cost, an out-going WFST arc and a current token)

```
1: procedure RECOMBINE(cost, arc, curTok)
2:     oldTokPack = state2tokPack[arc.next_state]
3:     curTokPack = packFunc(cost, arc.id)    ▷ pack into uint64
4:     ret = atomicMin (oldTokPack, curTokPack)
5:     if ret > curTokPack then               ▷ recombine
6:         perArcTokBuf[arc.id] = *curTok     ▷ store token
```

score=0.7
arc_id=1

score=0.8
arc_id=2

score=0.6
arc_id=3

score=0.9
arc_id=4

1. *atomicMax(*address, val)*. Computes *max(*address, val)*, writes the result to *address*, and returns the original *address* .

# 2nd Problem: Load balancing

- Load imbalance:  different num. of out-going arcs

# Load balancing

**Dispatcher**

idx=atomicAdd(global_d,1)



Group 1 with 32 threads
**Token ?**

t :ε

t(s):s

2

t(t):t

ε : t

7

(t)t:ε

(t)t(s):s

(t)s(s):s

(s)t(t):t

(t)s(t):t

6

(s)s(s):s

(s)t(s):s

(t)s:ε

(s)s(t):t

5

(s)t:ε

3

s(s):s

4

(s)s:ε

s(t):t

1

s :ε

# Load balancing

*Dispatcher*

idx=atomicAdd(global_d,1)

Group 1 with 32 threads
**Token ?**

t : ε

t(s):s

2

t(t):t

ε : t

7

(t)t:ε

(t)t(s):s

(t)s(s):s

(s)t(t):t

(t)s(t):t

6

(t)s:ε

(s)s(s):s

(s)t(s):s

(s)s(t):t

5

(s)t:ε

s(s):s

4

(s)s:ε

3

s(t):t

1

s : ε

# Load balancing

*Dispatcher*

idx=atomicAdd(global_d,1)

Group 1 with 32 threads
**Token 1 in State 1**

# Load balancing

# Load balancing



Group 0 with 32 threads
**Token ?**

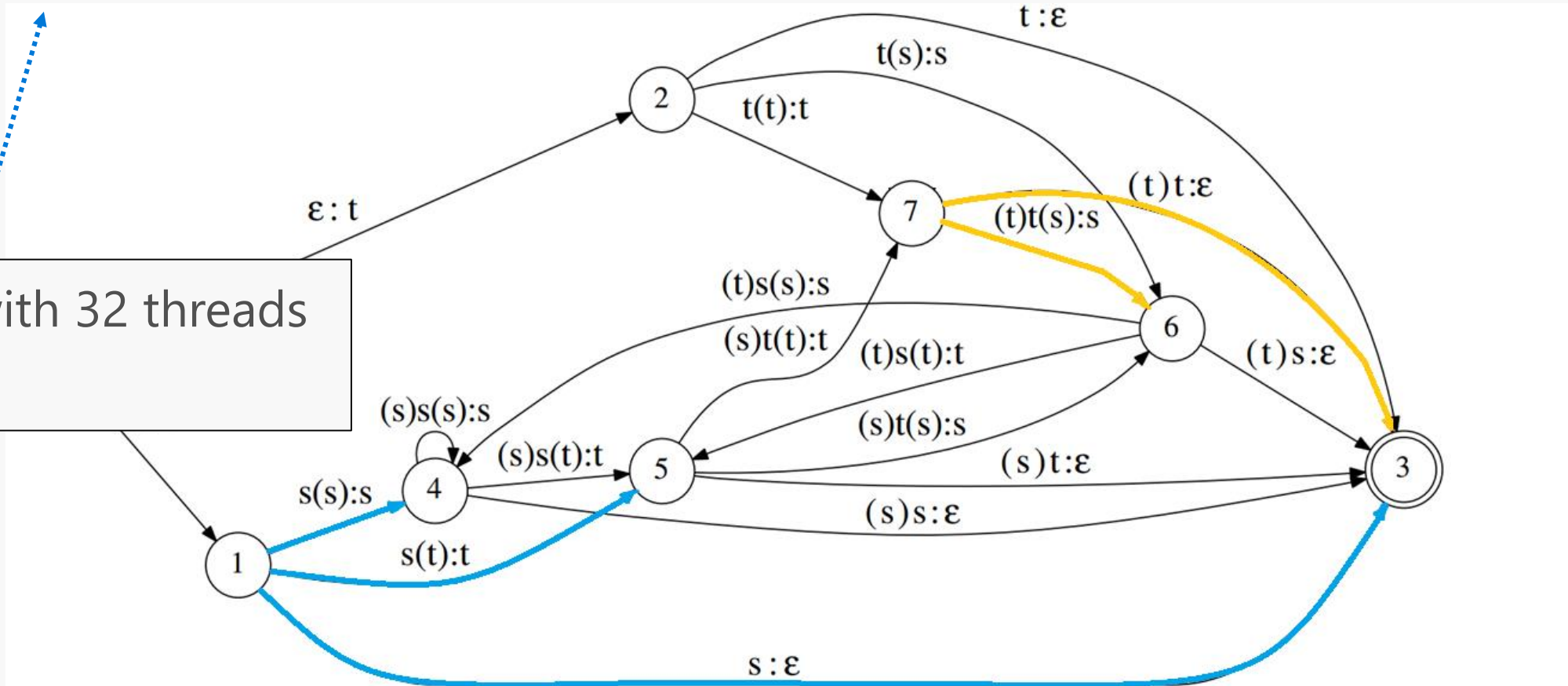Group 1 with 32 threads
**Token 1 in State 1**

# Load balancing



Group 0 with 32 threads
**Token ?**

Dispatcher
idx=atomicAdd(global_d,1)

Group 1 with 32 threads
**Token 1 in State 1**

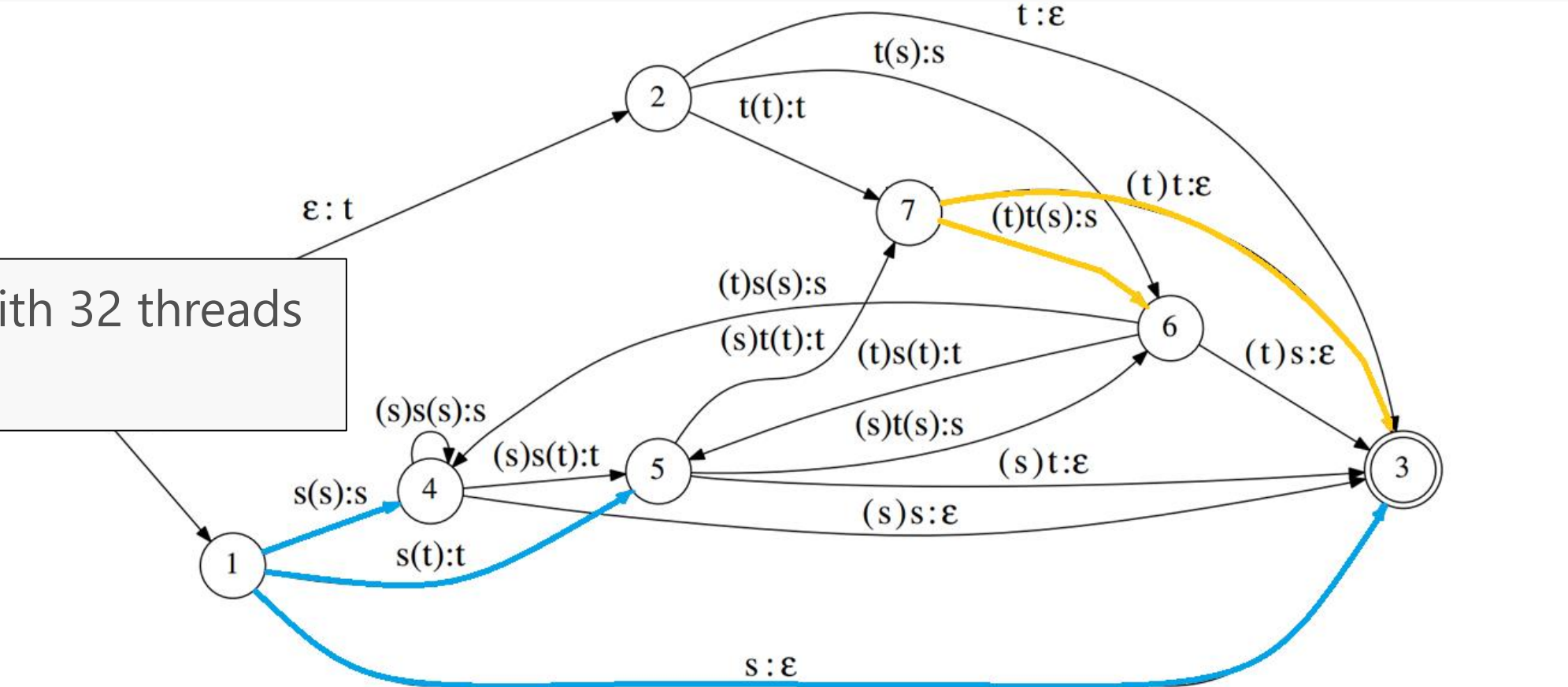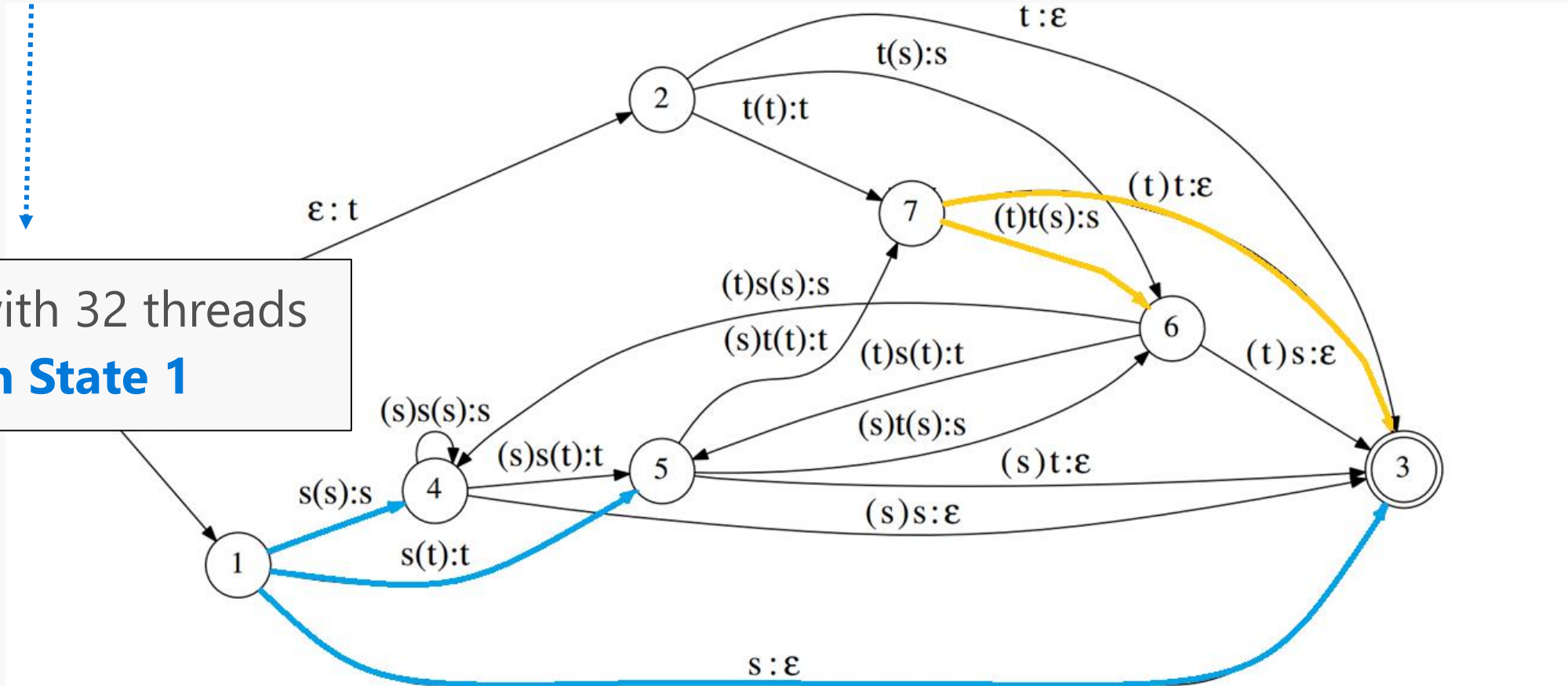# Load balancing



Group 0 with 32 threads
**Token 0 in State 2**

Group 1 with 32 threads
**Token ?**

# Load balancing

# Load balancing

- Dynamic load balancing



**Algorithm 2** Grid-level Token Passing (N=32; Inputs: the current active token vector)

1: **procedure** DYNAMIC LOAD BALANCING (toks)
2:     group = cooperative_groups::tiled_partition$\langle 32 \rangle$
3:     **if** group.thread_rank()==0 **then** ▷ rank 0 in each group
4:         i = *atomicAdd*(global_d,1) ▷ allocate new tokens
5:     i = group.*shfl*(i,0) ▷ rank 0 broadcasts i to whole group
6:     **if** i>=sizeof(toks) **then** return ▷ all tokens processed
7:     **for** arc in tok2arcs(toks[i]) **do** ▷ thread parallelism
8:         **call** *Recombine*(toks[i].cost+arc.cost, arc, toks[i])

# 3rd Problem: Lattice processing



N-best hypotheses:
1) reading time.
2) read in time.
3) re ding time.

- Linkedlist → vector
- Atomic operations
  - e.g. memory allocation

```
// implementation of v.push_back(val)
int idx = atomicAdd(cnt_d, 1); // idx=cnt_d++
mem_d[idx] = *val; // store data
```

- Parallel lattice pruning

# Experiments

- ## Setup
  - Switchboard 300 hours corpus, Cross Entropy & LF-MMI acoustic models (AM)
  - 30k-vocabulary, several tri-gram language models (LM)

- ## Baseline
  - Kaldi 1-best decoder
  - Kaldi lattice decoder

- ## GPU Optimization[1]
  - fast memcpy; merge GPU kernels by adding grid sync.; etc. (rel. 20% speedup)

1. https://github.com/chenzhehuai/kaldi/tree/gpu-decoder

# Experiments

- Performance

| system | lat. den. | WER | +rescored | OWER | NCE |
|--------|-----------|-----|-----------|------|-----|
| CPU | 30.3 | 15.5 | 14.3 | 11.2 | 0.322 |
| GPU | 30.2 | 15.5 | 14.3 | 11.2 | 0.328 |

Table 1: *1-best and Lattice Performance (beam=14).*

- The same 1-best & lattice quality

# Experiments

- Speedup

Table 2: *Speedup of the Proposed Method (beam=14).*

| system | 1-best | | + lattice | |
|---|---|---|---|---|
| | RTF | Δ | RTF | Δ |
| CPU | **0.16** | **1.0X** | **0.27** | **1.0X** |
| + 8-sequence (1 socket) | - | - | 0.15 | 1.8X |
| GPU | 0.016 | 10X | 0.080 | 3.3X |
| + atomic operation | 0.015 | 11X | 0.077 | 3.5X |
| + dyn. load balancing | **0.011** | **15X** | 0.075 | 3.6X |
| + lattice prune | - | - | **0.028** | **9.7X** |
| + 8-sequence (MPS) | 0.0035 | 46X | 0.0080 | 34X |

10X speedup

# Experiments

- Speedup

Table 2: *Speedup of the Proposed Method (beam=14).*

| system | 1-best | | + lattice | |
|---|---|---|---|---|
| | RTF | Δ | RTF | Δ |
| CPU | **0.16** | **1.0X** | **0.27** | **1.0X** |
| + 8-sequence (1 socket) | - | - | 0.15 | 1.8X |
| GPU | 0.016 | 10X | 0.080 | 3.3X |
| + atomic operation | 0.015 | 11X | 0.077 | 3.5X |
| + dyn. load balancing | **0.011** | **15X** | 0.075 | 3.6X |
| + lattice prune | - | - | **0.028** | **9.7X** |
| + 8-sequence (MPS) | 0.0035 | 46X | 0.0080 | 34X |

+ utterance-parallelism

# Experiments

- Speedup

Table 2: *Speedup of the Proposed Method (beam=14).*

| system | 1-best | | + lattice | |
|---|---|---|---|---|
| | RTF | Δ | RTF | Δ |
| CPU | **0.16** | **1.0X** | **0.27** | **1.0X** |
| + 8-sequence (1 socket) | - | - | 0.15 | 1.8X |
| GPU | 0.016 | 10X | 0.080 | 3.3X |
| + atomic operation | 0.015 | 11X | 0.077 | 3.5X |
| + dyn. load balancing | **0.011** | **15X** | 0.075 | 3.6X |
| + lattice prune | - | - | **0.028** | **9.7X** |
| + 8-sequence (MPS) | 0.0035 | 46X | 0.0080 | 34X |

+ utterance-parallelism

+ utterance-parallelism

# Experiments

- Speedup

Table 2: *Speedup of the Proposed Method (beam=14).*

| system | 1-best | | + lattice | |
|---|---|---|---|---|
| | RTF | Δ | RTF | Δ |
| CPU | **0.16** | **1.0X** | **0.27** | **1.0X** |
| + 8-sequence (1 socket) | - | - | 0.15 | 1.8X |
| GPU | 0.016 | 10X | 0.080 | 3.3X |
| + atomic operation | 0.015 | 11X | 0.077 | 3.5X |
| + dyn. load balancing | **0.011** | **15X** | 0.075 | 3.6X |
| + lattice prune | - | - | **0.028** | **9.7X** |
| + 8-sequence (MPS) | 0.0035 | 46X | 0.0080 | 34X |

**Improvement on naïve GPU decoder**

# Experiments

- Speedup

Table 2: *Speedup of the Proposed Method (beam=14).*

| system | 1-best | | + lattice | |
|---|---|---|---|---|
| | RTF | Δ | RTF | Δ |
| CPU | **0.16** | **1.0X** | **0.27** | **1.0X** |
| + 8-sequence (1 socket) | - | - | 0.15 | 1.8X |
| GPU | 0.016 | 10X | 0.080 | 3.3X |
| + atomic operation | 0.015 | 11X | 0.077 | 3.5X |
| + dyn. load balancing | **0.011** | **15X** | 0.075 | 3.6X |
| + lattice prune | - | - | **0.028** | **9.7X** |
| + 8-sequence (MPS) | 0.0035 | 46X | 0.0080 | 34X |

**Our new number[1] is over 50X**

1. https://github.com/chenzhehuai/kaldi/tree/gpu-decoder

# Experiments

- Compatibility
  - **GPU arch.**
  - **WFST size**
  - **Acoustic model**

# Experiments

- GPU Memory v.s. WFST size



**Enough for 1-pass WFST**

# Experiments

- Profiling



| | | |
|---|---|---|
| Others | D2H and H2D | Lattice pruning |
| Lattice generation | Token passing | |

# Conclusion & Future works

- Propose:
  - parallel Viterbi decoding & lattice processing

- Implementation:
  - Open-source & compatible with Kaldi recipes:
    https://github.com/chenzhehuai/kaldi/tree/gpu-decoder

- Future works:
  - More researches in GPU decoding
  - WFST algorithms, e.g. compose, determinize and minimize
  - Tight integration with acoustic inference (in GPUs)

# Backup materials

# Lattice processing

- ## Lattice pruning
  - The original CPU version: iteratively updates extra costs of nodes and arcs until they stop changing

  - Proposed:
    - updating in parallel
    - Linkedlist → vector
    - Atomic operations

**Algorithm 3** Grid-level Lattice Processing (processing frame, lattice token vector and lattice arc vector are taken as inputs)

1: **procedure** PRUNE LATTICE FOR FRAME (f, toks, arcs)
2:     **for** tok in toks(f-1) **do**    ▷ extraCost initialization
3:         tok.extraCost = FLT_MAX
4:     **while** modified == 1 **do**
5:         modified = 0
6:         **for** arc in arcs(f) **do**    ▷ thread parallelism
7:             cost = $ArcExtraCost$(arc)
8:                 ▷ returns the cost difference between the best path including the arc, and the best overall path.
9:             **if** cost < latticeBeam **then**
10:             $atomicMin$(tok.extraCost,cost)
11:             $atomicAdd$(modified,1)

# Speed and memory optimization

- share WFST between utterances in a GPU
- reduce context switching overheads: multi-process service (MPS)
- Lazy malloc to reduce memory: CudaMallocManaged
- grid sync implementation

```
__syncthreads();
if (threadIdx.x == 0) {
    int nb = 1;
    if (blockIdx.x == 0) {
        nb = 0x80000000 - (gridDim.x-1);
    }
    int old_epoch = *fast_epoch;
    __threadfence();
    atomicAdd((int*)fast_epoch, nb);
    int cnt=0;
    while (((*fast_epoch) ^ old_epoch) >= 0) ;
}
__syncthreads();
```

# Speed and memory optimization

- reduce grid sync using multiple copy of variable
  - Do not need to make threads wait for modified1 = false

```
do {
    __grid_sync_nv_internal(params.barrier);

    swap(modified0,modified1);
    *modified1 = false;

    //__grid_sync_nv_internal(params.barrier);

    processNonEmittingTokens_function<32,2>(params,cutoff,size,modified0);
    __grid_sync_nv_internal(params.barrier);
} while ((*modified0)==true);
```

# Speed and memory optimization

- reduce grid sync using multiple copy of variable
- faster atomicAdd by launch-and-go

```
//idx=atomicAdd((int*)fast_epoch, nb);
atomicAdd((int*)fast_epoch, nb);
```

- sharing atomicAdd
- fast copy

```
asm("st.global.v2.u64 [%0], {%1,%2};" :: "l"(a), "l"(src.x), "l"(src.y));
```