Multi-Threaded Composition of Weighted Finite-State Transducers

Bryan Jurish

Berlin-Brandenburg Academy of Sciences

jurish@bbaw.de

Kay-Michael Würzner

University of Potsdam wuerzner@uni-potsdam.de

WATA 2012 Dresden, 29th May, 2012



berlin-brandenburgische AKADEMIE DER WISSENSCHAFTEN Deutsche Forschungsgemeinschaft

WATA 2012 / Jurish | Würzner / Multi-threaded composition of WFSTs - p. 1/24



The Big Idea

- The Situation
- The Approach

Parallel Composition Algorithms

- Master-Slave
- Peer-to-Peer

Experiments

- Materials
- Method
- Results

Concluding Remarks





— The Big Idea —





Deutsche Forschungsgemeinschaft

WATA 2012 / Jurish | Würzner / Multi-threaded composition of WFSTs - p. 3/24

The Situation

No Free Lunch (anymore)

- CPU frequency growth stagnating
- Multiprocessor systems increasingly popular
- → *"horizontal" scaling* / multi-threading

(W)FST Composition

 $\mathcal{T}_3 = (\mathcal{T}_1 \circ \mathcal{T}_2)$

- Online: lexical lookup, Viterbi decoding, parsing, ...
- Offline: lexicon compilation, statistical modelling, ...
- no generic parallel implementation (that we know of)

Amdahl's Law

$$S(N) = \frac{1}{(1-P) + \frac{P}{N}}$$

- Not all algorithms scale well horizontally ($P \ll 1$)
- For WFSTs, *P* may depend on *FST topology not all WFST compositions scale horizontally!*





Deutsche Forschungsgemeinschaft

DEG

The Basics

Definition

Given two ε -free WFSTs $\mathcal{T}_1 = \langle \Sigma, \Gamma, Q_1, q_{0_1}, F_1, E_1, \rho_1 \rangle$ and $\mathcal{T}_2 = \langle \Gamma, \Delta, Q_2, q_{0_2}, F_2, E_2, \rho_2 \rangle$, $\mathcal{T}_3 = (\mathcal{T}_1 \circ \mathcal{T}_2)$ is itself a WFST, with:

$$\begin{bmatrix} \mathcal{T}_3 \end{bmatrix} : (x, y) \mapsto \bigoplus_{z \in \Gamma^*} \llbracket \mathcal{T}_1 \rrbracket (x, z) \otimes \llbracket \mathcal{T}_2 \rrbracket (z, y) \\ \mathcal{T}_3 = \langle \Sigma, \Delta, Q_1 \times Q_2, (q_{0_1}, q_{0_2}), F_1 \times F_2, E_3, \rho_1 \otimes \rho_2 \rangle \\ E_3 = \biguplus_{\substack{(q_1, r_1, a, b, w_1) \in E_1, \\ (q_2, r_2, b, c, w_2) \in E_2}} \left\{ ((q_1, q_2), (r_1, r_2), a, c, w_1 \otimes w_2) \right\}$$

Properties

- simple construction requires ε-free WFSTs
- commutative & complete semiring $\langle \mathbb{K}, \oplus, \otimes, \overline{1}, \overline{0} \rangle$
- worst-case $\mathcal{O}_{\text{time}} = \mathcal{O}(|E_1 \times E_2|)$





Serial Algorithm

 $\texttt{compose}\big(\mathcal{T}_1 = \langle \Sigma, \Gamma, Q_1, q_{0_1}, F_1, E_1, \rho_1 \rangle, \mathcal{T}_2 = \langle \Gamma, \Delta, Q_2, q_{0_2}, F_2, E_2, \rho_2 \rangle\big)$ **1** $Q \leftarrow \{(q_{0_1}, q_{0_2})\}$ /* initialize */ **2** $V \leftarrow \{(q_{0_1}, q_{0_2})\}$ /* visitation queue */ 3 while $V \neq \emptyset$ do $(q_1, q_2) \leftarrow \texttt{pop}(V)$ 4 /* visit state */ if $(q_1,q_2)\in F_1 imes F_2$ then 5 /* final state */ $F \leftarrow F \cup \{(q_1, q_2)\}$ 6 $\left| \rho(q_1, q_2) \leftarrow \rho_1(q_1) \otimes \rho_2(q_2) \right|$ 7 foreach $(e_1, e_2) \in E[q_1] \times E[q_2]$ with $o[e_1] = i[e_2]$ do 8 /* align edges */ if $(n[e_1], n[e_2]) \notin Q$ then 9 $Q \leftarrow Q \cup \{(\mathbf{n}[e_1], \mathbf{n}[e_2])\}$ 10 $V \leftarrow V \cup \{(\mathbf{n}[e_1], \mathbf{n}[e_2])\}$ 11 /* enqueue for visitation */ $E \leftarrow E \cup \{(q_1, q_2), (n[e_1], n[e_2]), i[e_1], o[e_2], w[e_1] \otimes w[e_2]\}$ 12 13 return $\mathcal{T}_3 = \langle \Sigma, \Delta, Q, (q_{0_1}, q_{0_2}), F, E, \rho \rangle$





The Approach

Parallel State Visitation

- breadth-first search of output states
- distributed output data
- shared visitation queue

(lines 4–12)
(
$$V$$
 : FIFO)
(Q, F, E, ρ)
(V)

Amdahl's Law Revisited



assumes constant (average) state complexity

Deutsche

DFG

Forschungsgemeinschaft

worst-case breadth-first visitation



berlin-brandenburgi

AKADEMIE DER WISSENSCHAFTE



WATA 2012 / Jurish | Würzner / Multi-threaded composition of WFSTs - p. 7/24

- Algorithms -





Deutsche Forschungsgemeinschaft

WATA 2012 / Jurish | Würzner / Multi-threaded composition of WFSTs - p. 8/24

Algorithm (Sketch): Master-Slave



Superordinate Distribution of Work

state-pairs (q_1, q_2) passed to slaves for visitation

Slave Tasks

align & expand transitions, globally enqueue visitation requests

Shared Global Data

V : visitation queue

berlin-brandenburgis

AKADEMIE DER WISSENSCHAFTER

- \mathbf{Q} : visited states
- n_q : output state counter
- n_up: number of tasks currently assigned



Deutsche Forschungsgemeinschaft

Algorithm (Sketch): Peer-to-Peer



State Partitioning Function

• peer *i* visits states with $r(q_1, q_2) = i$

Peer-to-Peer Message Passing

- messages are aligned transitions (e_1, e_2)
- sender: $r(p[e_1], p[e_2]) \rightarrow receiver: r(n[e_1], n[e_2])$

Shared Global Data

- n_q : output state counter
- n_up: number of messages currently enqueued (for termination)

(for serialization)

 $V \in \wp(E_1 \times E_2)^{N \times N}$

 $r: (q_1, q_2) \mapsto \left| \frac{q_1 + q_2}{2} \right| \mod N$



AKADEMIE DER WISSENSCHAFTI

Deutsche Forschungsgemeinschaft

- Experiments -





Deutsche Forschungsgemeinschaft

WATA 2012 / Jurish | Würzner / Multi-threaded composition of WFSTs - p. 11/24

Experiments

Materials

- 2,266 randomly generated WFSTs \mathcal{T}
 - trie spine + random arcs
 - (piecewise-) uniform sampling
 - "embarrassingly parallel" topology
- algorithms implemented in C++
- hexadecacore test machine

```
\begin{aligned} \operatorname{depth}(\mathcal{T}) &\leq 32 \\ |Q_{\mathcal{T}}|, |E_{\mathcal{T}}|, |\Sigma| \\ P_{(\mathcal{T}^{-1} \circ \mathcal{T})} &> 99\% \end{aligned}
```

g++ v4.4.5 16 physical cores

Method

- for each generated \mathcal{T} , compute $(\mathcal{T}^{-1} \circ \mathcal{T})$
 - sample selection filter
 - varied number of threads

Evaluation

- average running time
- structural properties of $\mathcal{T}, (\mathcal{T}^{-1} \circ \mathcal{T})$

DFG

 $\frac{1}{64} \sec \le t_{\text{serial}} \le 8 \sec N \in \{1, 2, 4, 8, 16\}$



Deutsche Forschungsgemeinschaft $|Q|, |E|, \ldots$

8 iterations per configuration

Results: Master-Slave



S = t.serial / t.ms

Sourcestrain

Results: Peer-to-Peer



Concluding Remarks

Summary

- No (more) Free Lunch
 - parallelization of "traditional" serial algorithms
- Amdahl's Law Applied
 maximum speedup depends on FST topology
- Sharing (data) Hurts
 - distributed synchronization improves performance

Future Directions

- improve sampling procedure
- reduce p2p communication overhead?
- massive parallel architectures (CUDA, OpenCL)?



berlin-brandenburgische AKADEMIE DER WISSENSCHAFTEN

The End

Thank you for listening!





Deutsche Forschungsgemeinschaft

WATA 2012 / Jurish Würzner / Multi-threaded composition of WFSTs - p. 16/24



2d Plots

- $t_{\text{serial}}: S$
- E:S
- E/Q: S

3d Plots

- E/Q: N: S
- $t_{\text{serial}}: N: S$
- Q: E: histogram
- $t_{\text{serial}}: E/Q:$ histogram



berlin-brandenburgische AKADEMIE DER WISSENSCHAFTEN



Plots: 2d: $t_{serial} : S$



miversita.

Plots: 2d: *E* : *S*



universita.

Plots: 2d: E/Q : S



Plots: 3d: E/Q : N : S

master-slave

peer-to-peer





berlin-brandenburgische AKADEMIE DER WISSENSCHAFTEN

Plots: 3d: $t_{\text{serial}} : N : S$

master-slave

peer-to-peer





berlin-brandenburgische AKADEMIE DER WISSENSCHAFTEN

Plots: 3d: Q : E : **histogram**







raw

Deutsche Forschungsgemeinschaft

WATA 2012 / Jurish Würzner / Multi-threaded composition of WFSTs - p. 23/24

smoothed

Plots: 3d: t_{serial} : E/Q : histogram

raw

smoothed





berlin-brandenburgische AKADEMIE DER WISSENSCHAFTEN

Deutsche Forschungsgemeinschaft

WATA 2012 / Jurish | Würzner / Multi-threaded composition of WFSTs - p. 24/24