

Multi-Threaded Composition of Finite-State Automata

Bryan Jurish

Berlin-Brandenburg Academy of Sciences

jurish@bbaw.de

Kay-Michael Würzner

University of Potsdam

wuerzner@uni-potsdam.de

FSMNLP 2013

St. Andrews, 17th July, 2013

Overview

The Big Idea

- The Situation
- The Approach

Parallel Composition Algorithms

- Master-Slave
- Peer-to-Peer

Experiments

- Materials
- Method
- Results

Concluding Remarks

— The Big Idea —

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 FSTs, P may depend on *FST topology*
- ↪ *not all FST compositions scale horizontally!*

The Basics

Definition

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

$$\begin{aligned}\mathcal{T}_3 &= \langle \Sigma, \Delta, Q_1 \times Q_2, (q_{0_1}, q_{0_2}), F_1 \times F_2, E_3 \rangle \\ E_3 &= \bigcup_{\substack{(q_1, r_1, a, b) \in E_1, \\ (q_2, r_2, b, c) \in E_2}} \left\{ ((q_1, q_2), (r_1, r_2), a, c) \right\} \\ \llbracket \mathcal{T}_3 \rrbracket &= \{ (x, z) \mid \exists y : (x, y) \in \llbracket \mathcal{T}_1 \rrbracket \ \& \ (y, z) \in \llbracket \mathcal{T}_2 \rrbracket \} \\ &= \llbracket \mathcal{T}_1 \rrbracket \circ \llbracket \mathcal{T}_2 \rrbracket\end{aligned}$$

Properties

- simple construction requires ε -free FSTs
- worst-case $\mathcal{O}_{\text{time}} = \mathcal{O}(|E_1 \times E_2|)$

Serial Algorithm

```
compose( $\mathcal{T}_1 = \langle \Sigma, \Gamma, Q_1, q_{0_1}, F_1, E_1 \rangle, \mathcal{T}_2 = \langle \Gamma, \Delta, Q_2, q_{0_2}, F_2, E_2 \rangle$ )
```

```
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
```

```
4    $(q_1, q_2) \leftarrow \text{pop}(V)$  /* visit state */
```

```
5   if  $(q_1, q_2) \in F_1 \times F_2$  then /* final state */
```

```
6      $F \leftarrow F \cup \{(q_1, q_2)\}$ 
```

```
7     foreach  $(e_1, e_2) \in E[q_1] \times E[q_2]$  with  $o[e_1] = i[e_2]$  do /* align edges */
```

```
8       if  $(n[e_1], n[e_2]) \notin Q$  then
```

```
9          $Q \leftarrow Q \cup \{(n[e_1], n[e_2])\}$ 
```

```
10         $V \leftarrow V \cup \{(n[e_1], n[e_2])\}$  /* enqueue for visitation */
```

```
11         $E \leftarrow E \cup \{(q_1, q_2), (n[e_1], n[e_2]), i[e_1], o[e_2]\}$ 
```

```
12 return  $\mathcal{T}_3 = \langle \Sigma, \Delta, Q, (q_{0_1}, q_{0_2}), F, E \rangle$ 
```

The Approach

Parallel State Visitation

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

(lines 4–11)

(V : FIFO)

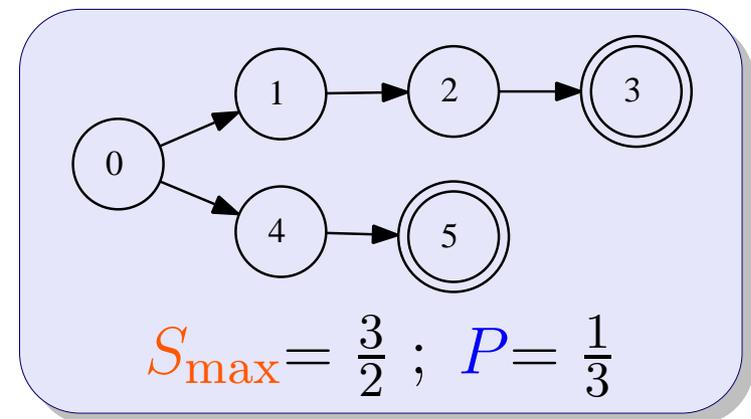
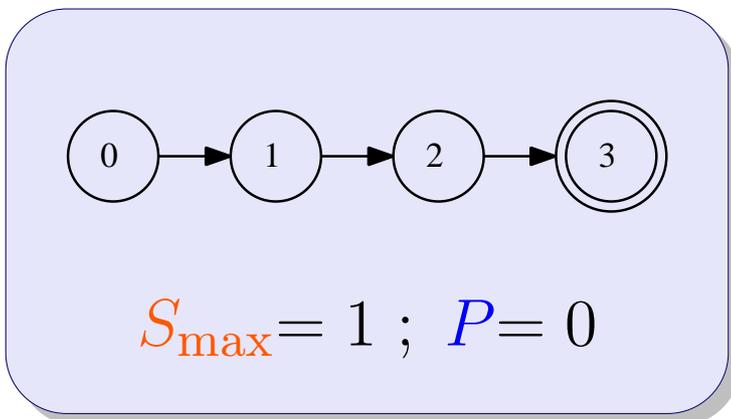
(Q, F, E)

(V)

Amdahl's Law Revisited

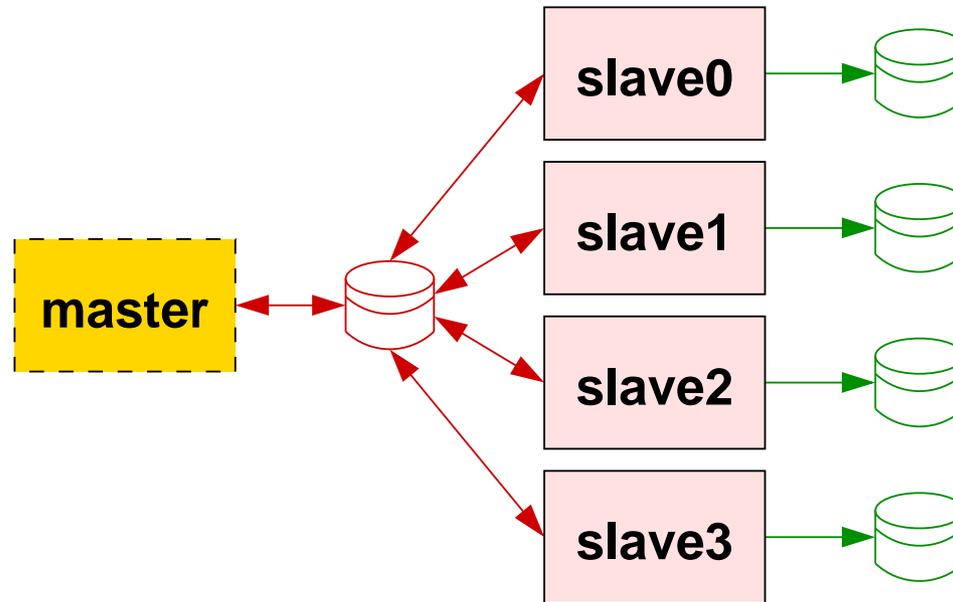
$$S_{\max} \approx \frac{|Q|}{1 + \text{depth}(\mathcal{T}_3)} = \frac{|Q|}{1 + \max_{q \in Q} \min_{\pi \in \Pi(q_0, q)} |\pi|}$$
$$P = 1 - \frac{1}{S_{\max}}$$

- assumes constant (average) state complexity
- worst-case breadth-first visitation



— Algorithms —

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
- Q : visited states
- n_q : output state counter
- n_{up} : number of tasks currently assigned

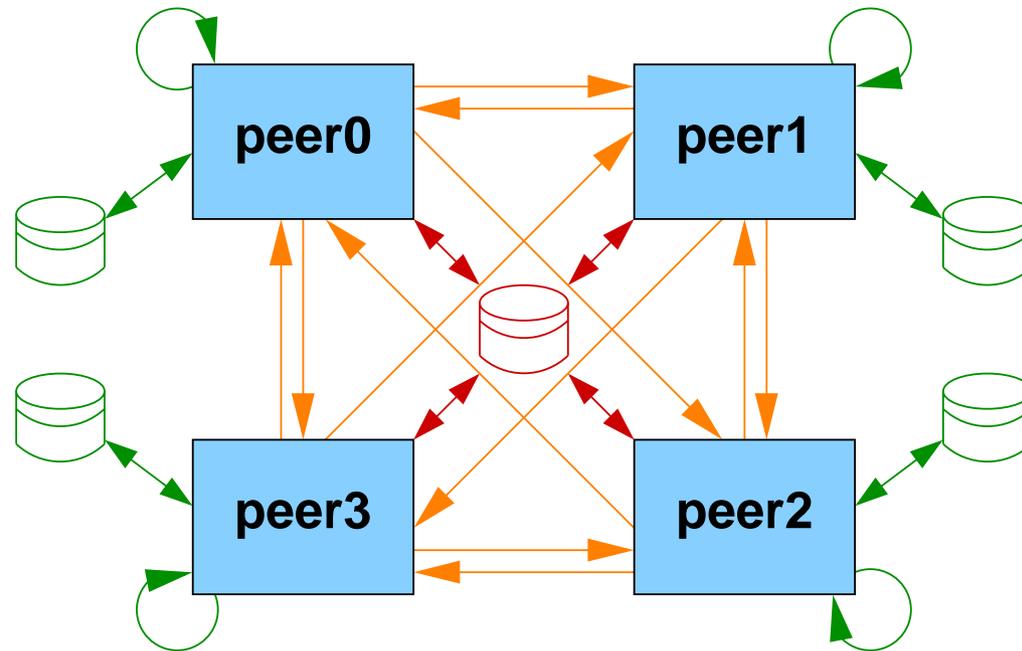
$$V \subseteq Q_1 \times Q_2$$

$$Q \subseteq Q_1 \times Q_2$$

(for serialization)

(for termination)

Algorithm (Sketch): Peer-to-Peer



State Partitioning Function

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

$$r : (q_1, q_2) \mapsto \left\lfloor \frac{q_1 + q_2}{2} \right\rfloor \bmod N$$

Peer-to-Peer Message Passing

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

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

Shared Global Data

- n_q : output state counter *(for serialization)*
- n_{up} : number of messages currently enqueued *(for termination)*

— Experiments —

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
- $\text{depth}(\mathcal{T}) \leq 32$
 $|Q_{\mathcal{T}}|, |E_{\mathcal{T}}|, |\Sigma|$
 $P_{(\mathcal{T}^{-1} \circ \mathcal{T})} > 99\%$
g++ v4.4.5
16 hardware cores

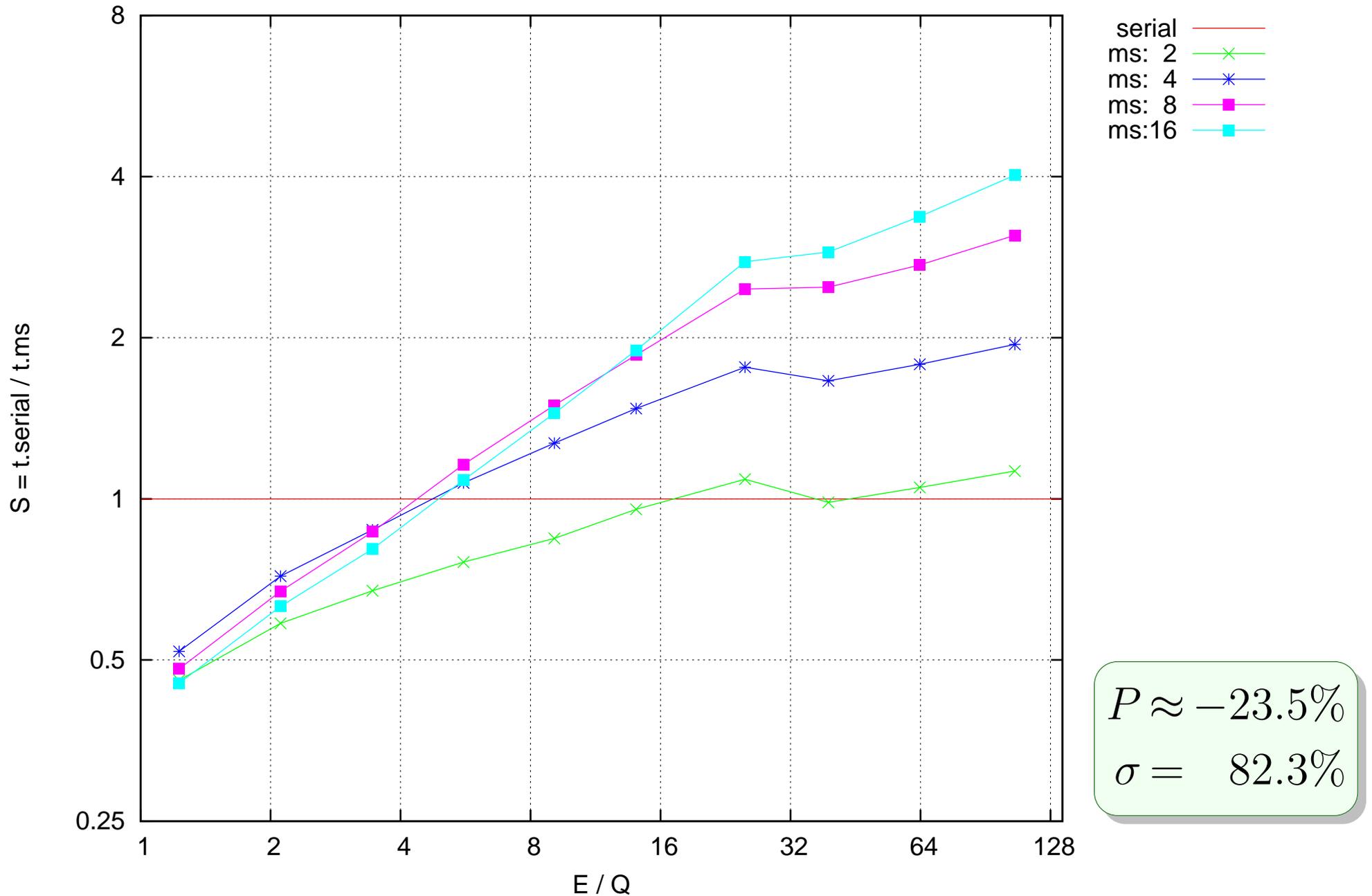
Method

- for each generated \mathcal{T} , compute $(\mathcal{T}^{-1} \circ \mathcal{T})$
 - ▶ sample selection filter
 - ▶ varied number of threads
- $\frac{1}{64} \text{ sec} \leq t_{\text{serial}} \leq 8 \text{ sec}$
 $N \in \{1, 2, 4, 8, 16\}$

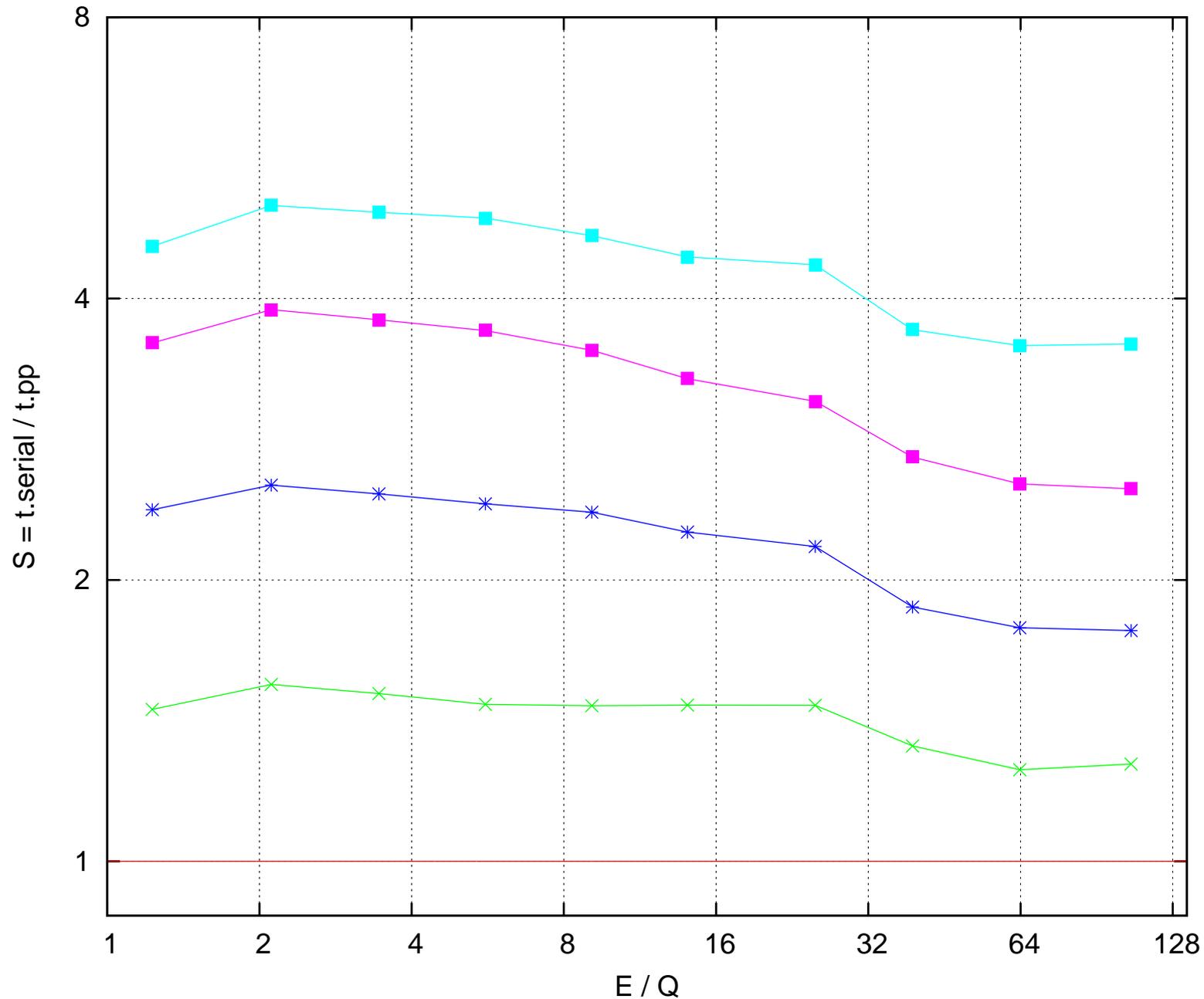
Evaluation

- average running time *8 iterations per configuration*
- structural properties of $\mathcal{T}, (\mathcal{T}^{-1} \circ \mathcal{T})$ $|Q|, |E|, \dots$

Results: Master-Slave



Results: Peer-to-Peer



$P \approx 83.1 \%$
 $\sigma = 7.18\%$

So What About NLP?

Lexical Lookup

- many “small” compositions
- topology-dependent S_{\max}

$$(\text{Id}(w) \circ \mathcal{T}_{\text{Lex}})_{w \in W}$$

\rightsquigarrow *prefer high-level fork() over W*

Corpus Analysis

- single “large” composition
- distributed representation

$$\mathcal{A}_{\text{Corpus}} \circ \mathcal{T}_{\text{Anal}}$$

\rightsquigarrow *serialization overhead*

Model Compilation

- offline “large” composition
- partitioning function

$$\mathcal{T}_{\text{Error}} \circ \mathcal{A}_{\text{Lex}}$$

\rightsquigarrow *task-dependent tuning*

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
- extend to other FST operations
 - ▶ determinization
 - ▶ minimization
 - ▶ cascaded best-path lookup

The End

Thank you for listening!

— Addenda —

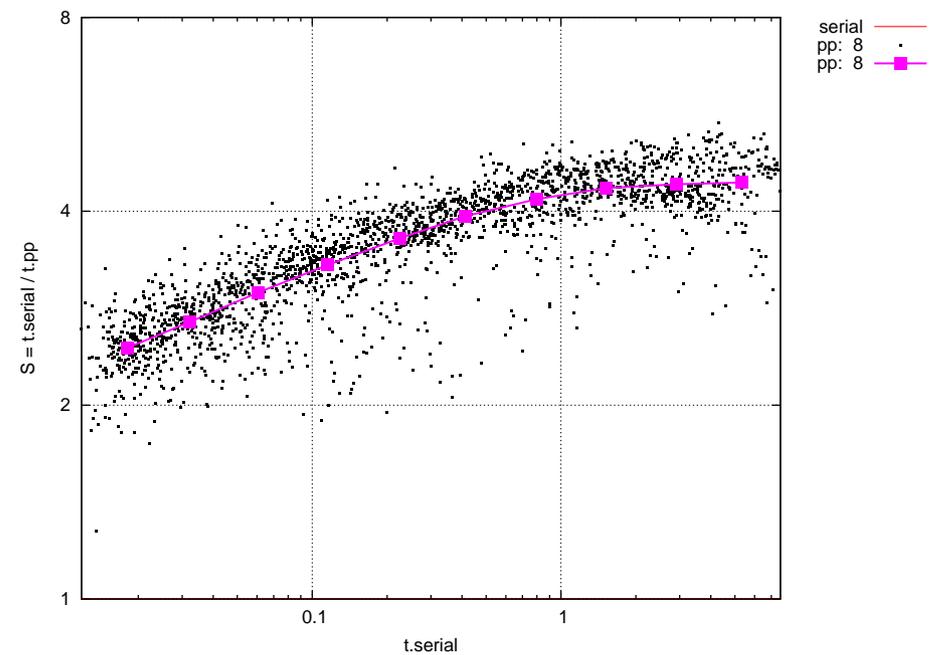
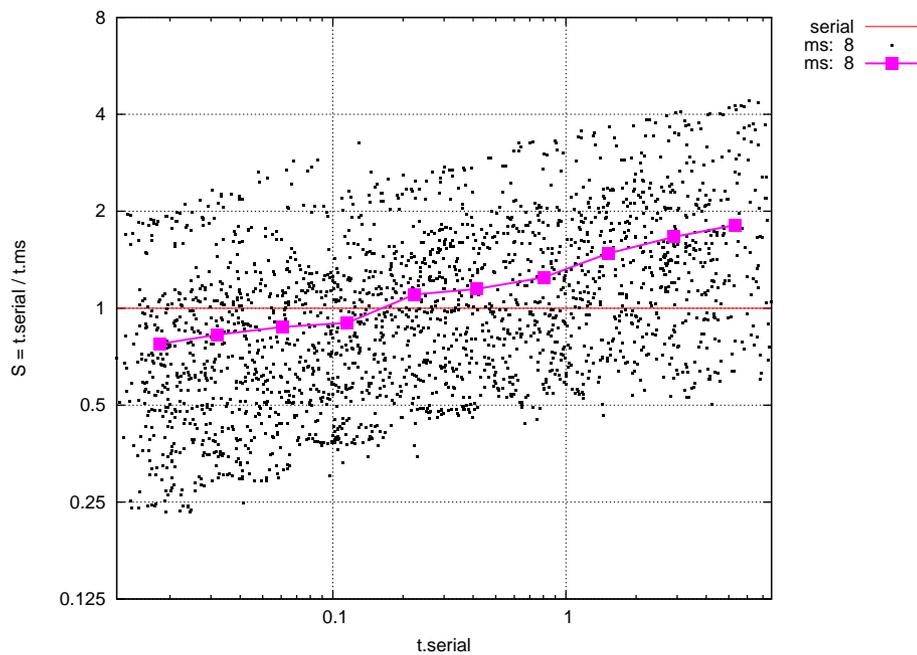
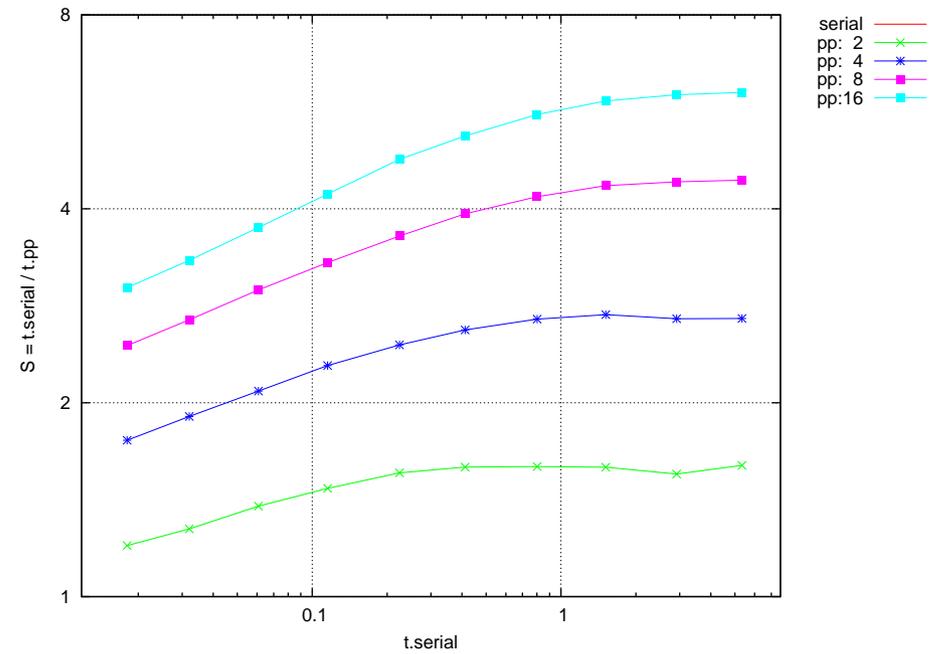
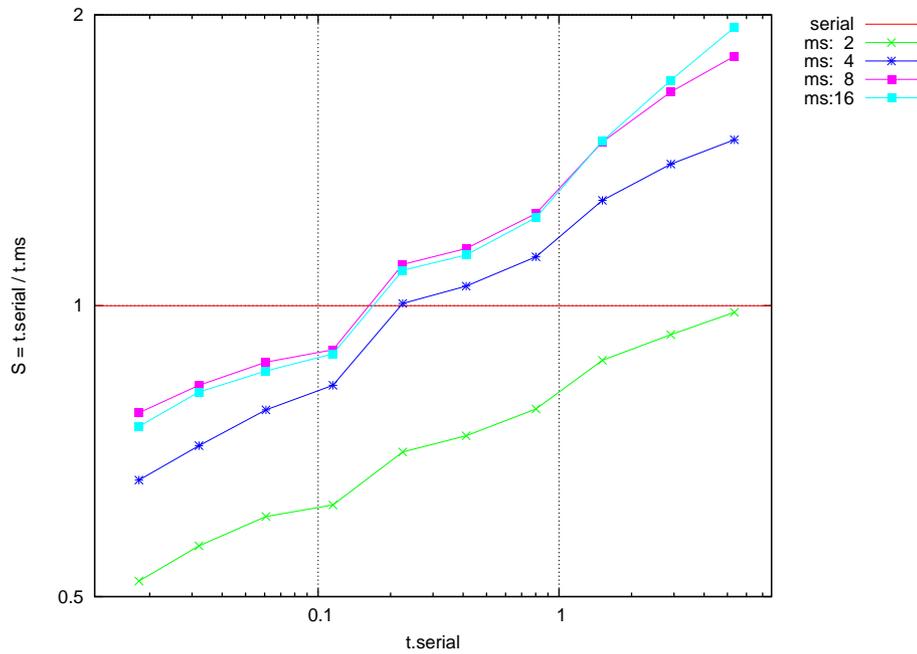
2d Plots

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

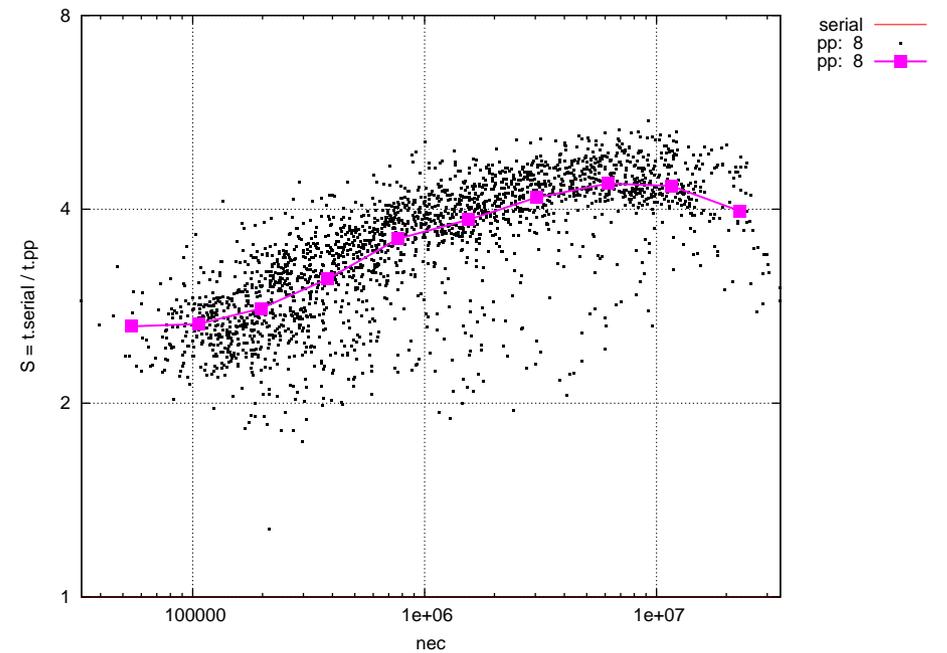
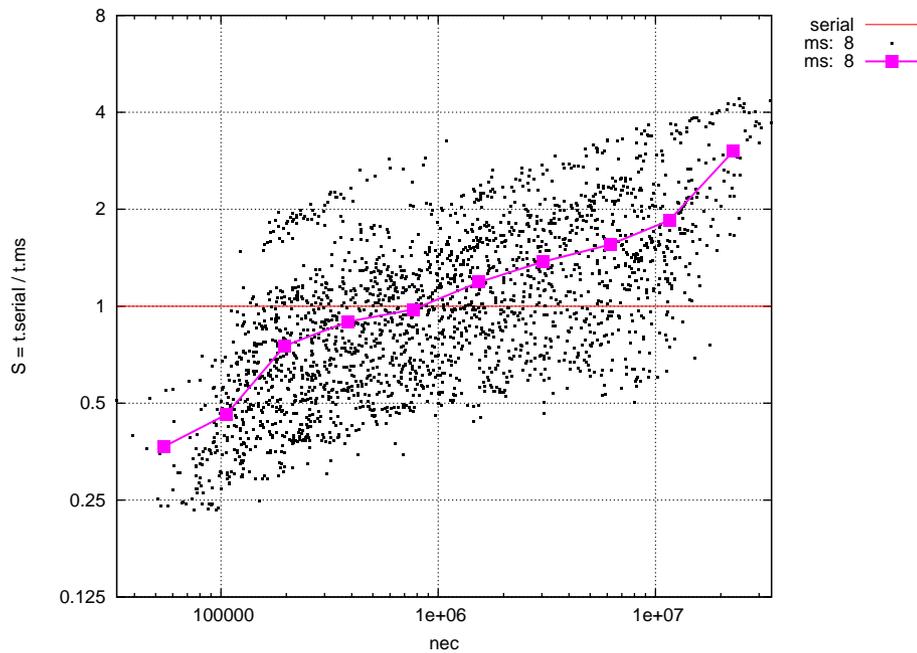
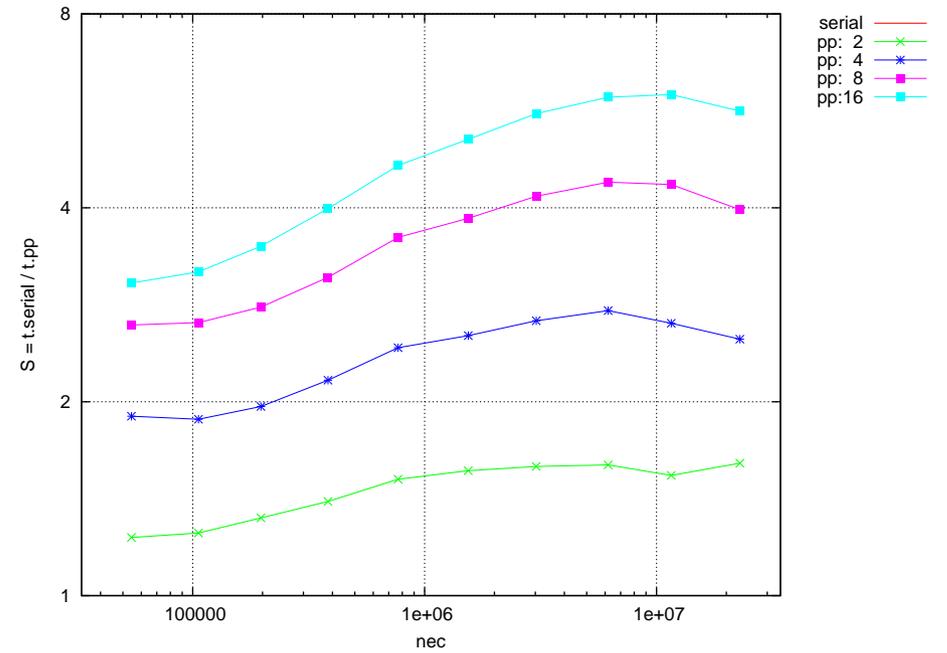
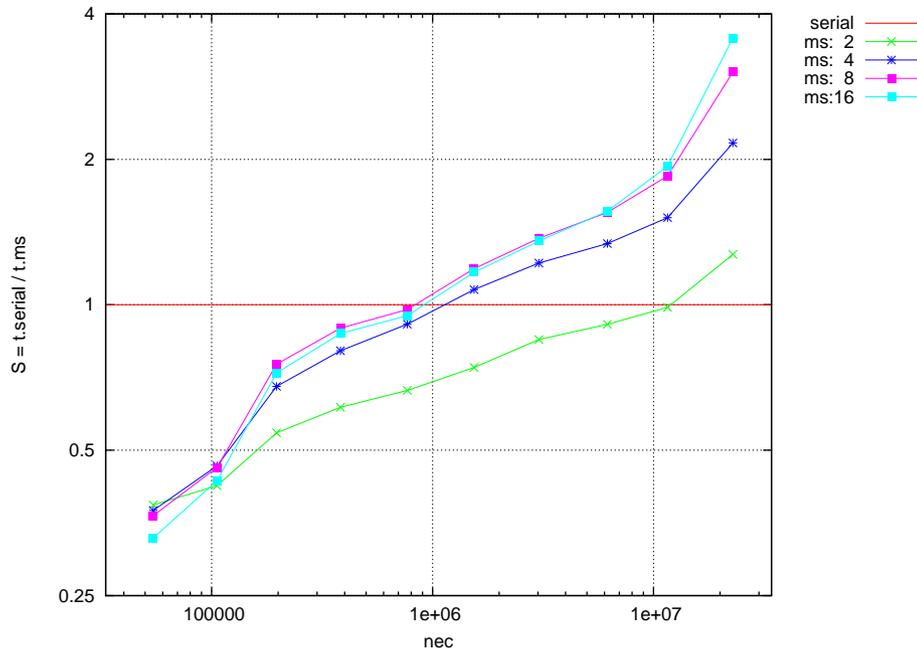
3d Plots

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

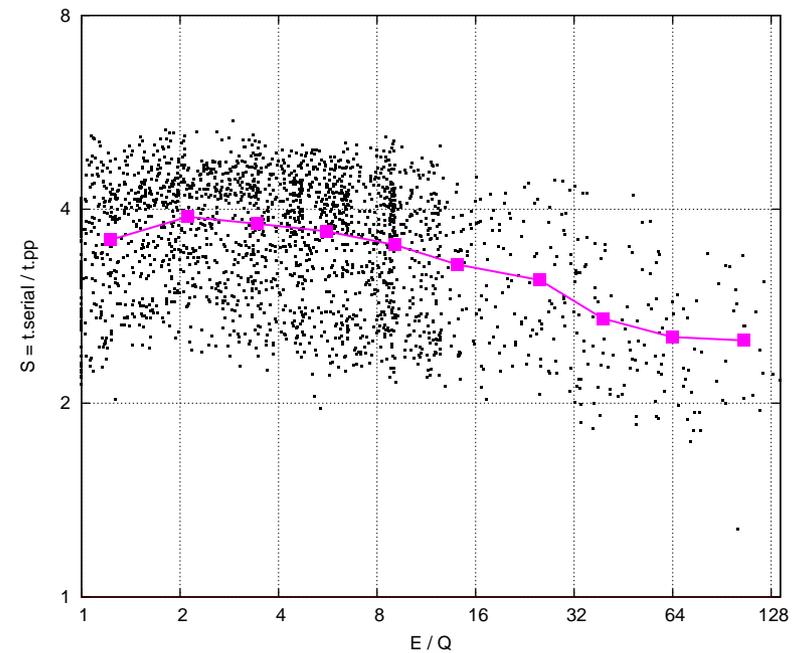
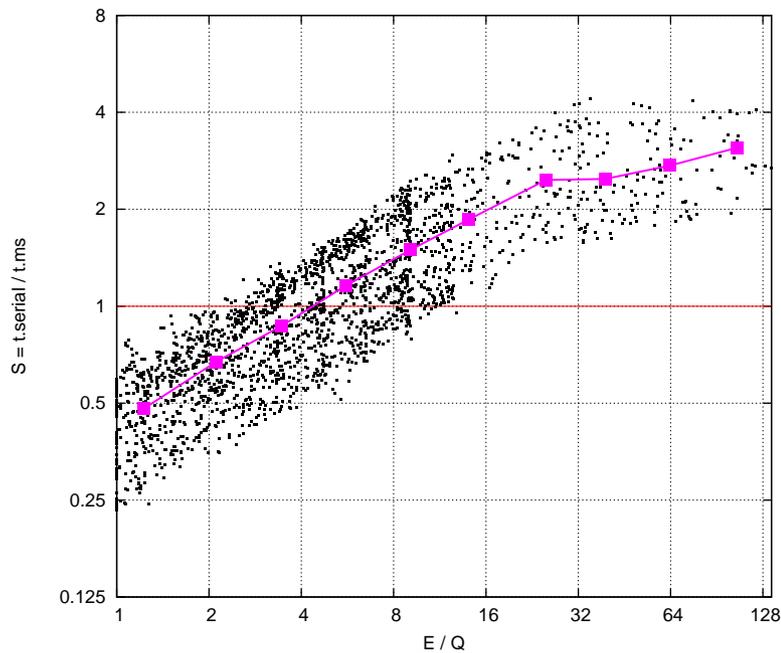
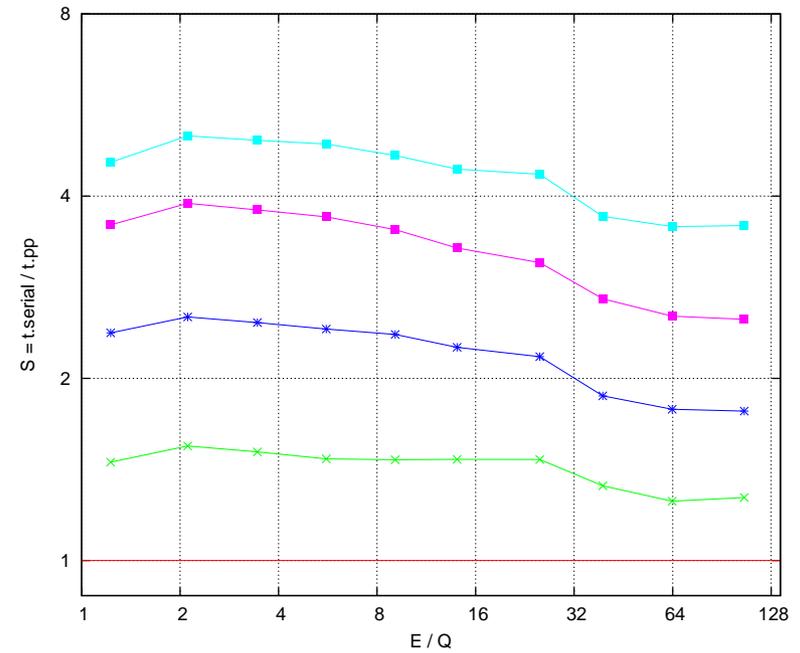
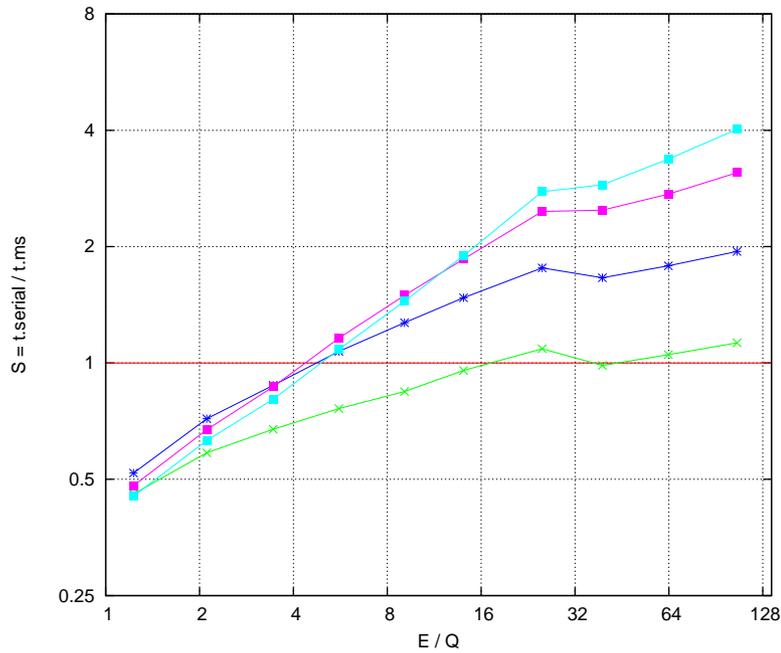
Plots: 2d: $t_{\text{serial}} : S$



Plots: 2d: $E : S$



Plots: 2d: $E/Q : S$



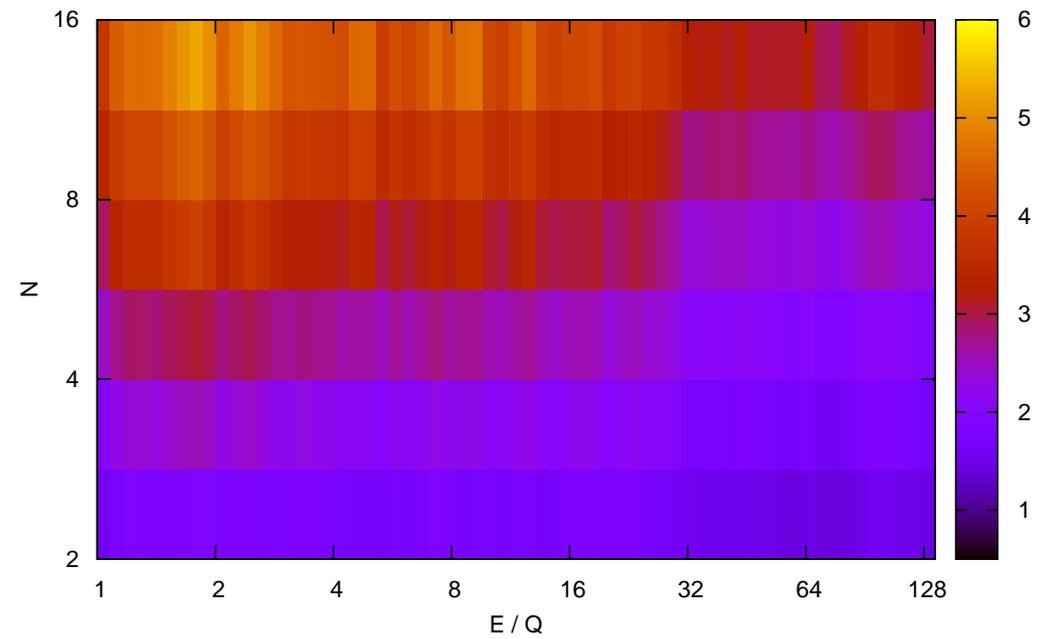
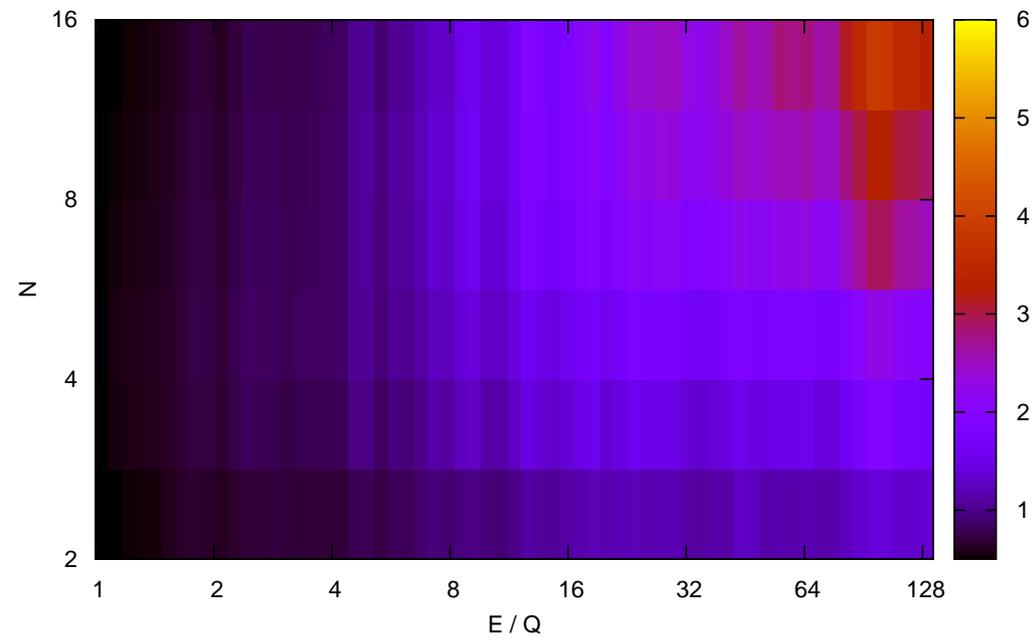
Plots: 3d: $E/Q : N : S$

master-slave

peer-to-peer

$S = t.\text{serial} / t.\text{ms}$

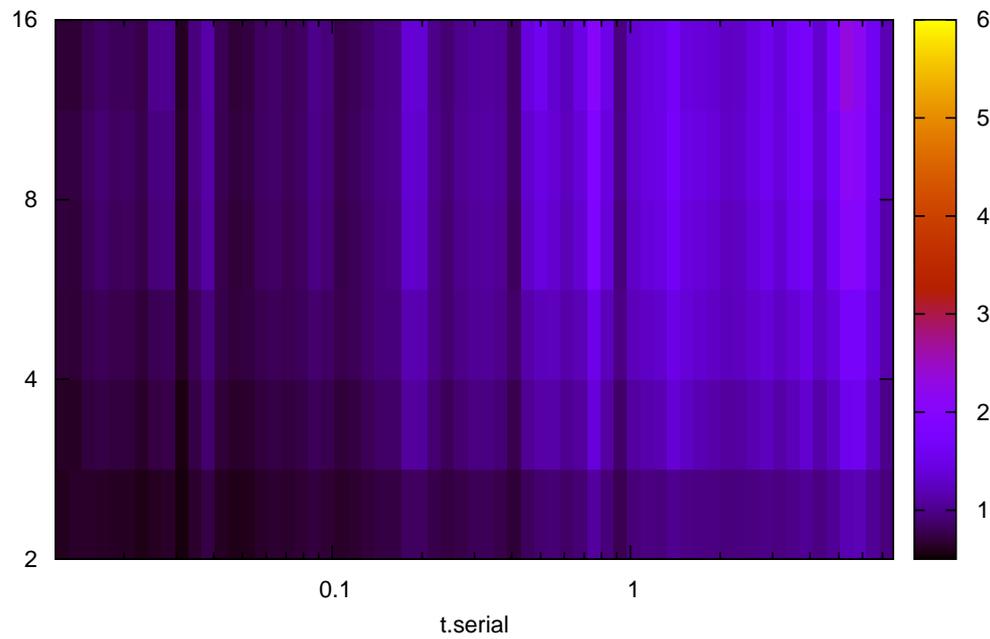
$S = t.\text{serial} / t.\text{pp}$



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

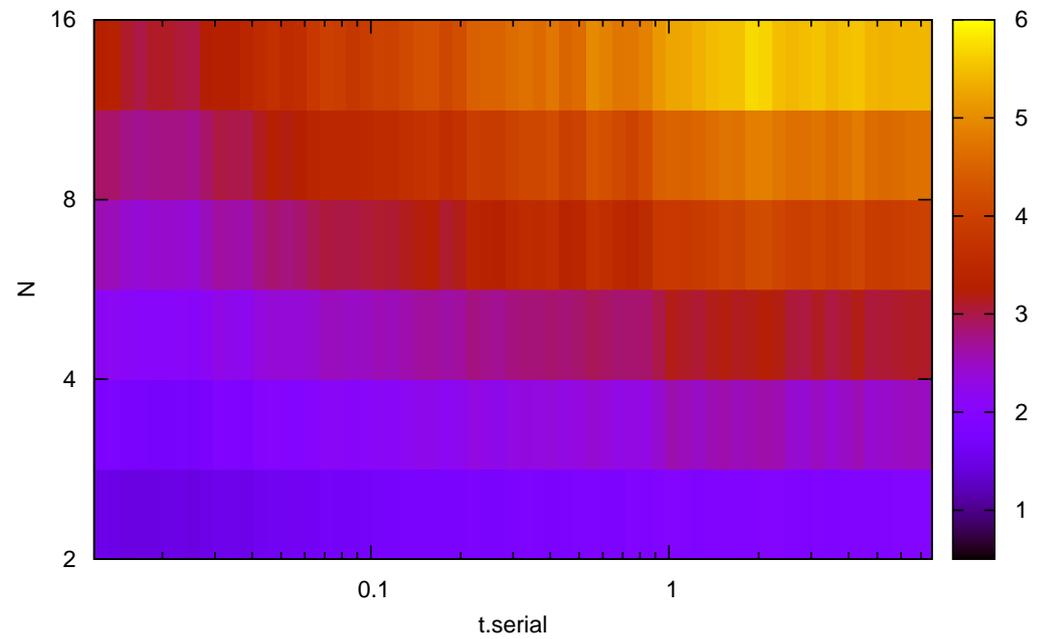
master-slave

$S = t_{\text{serial}} / t_{\text{ms}}$



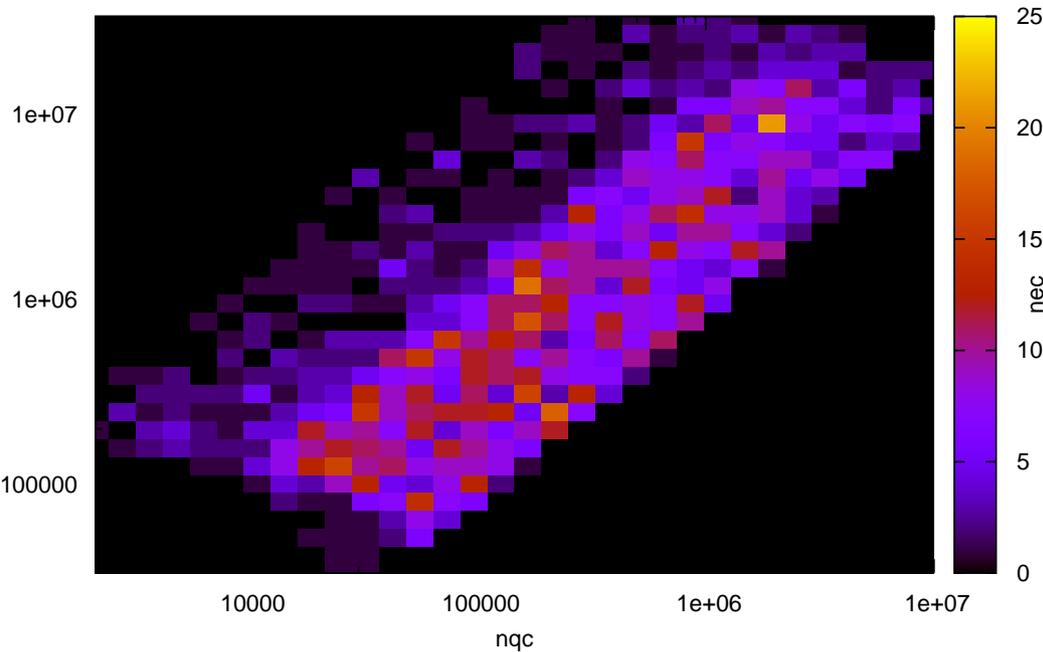
peer-to-peer

$S = t_{\text{serial}} / t_{\text{pp}}$

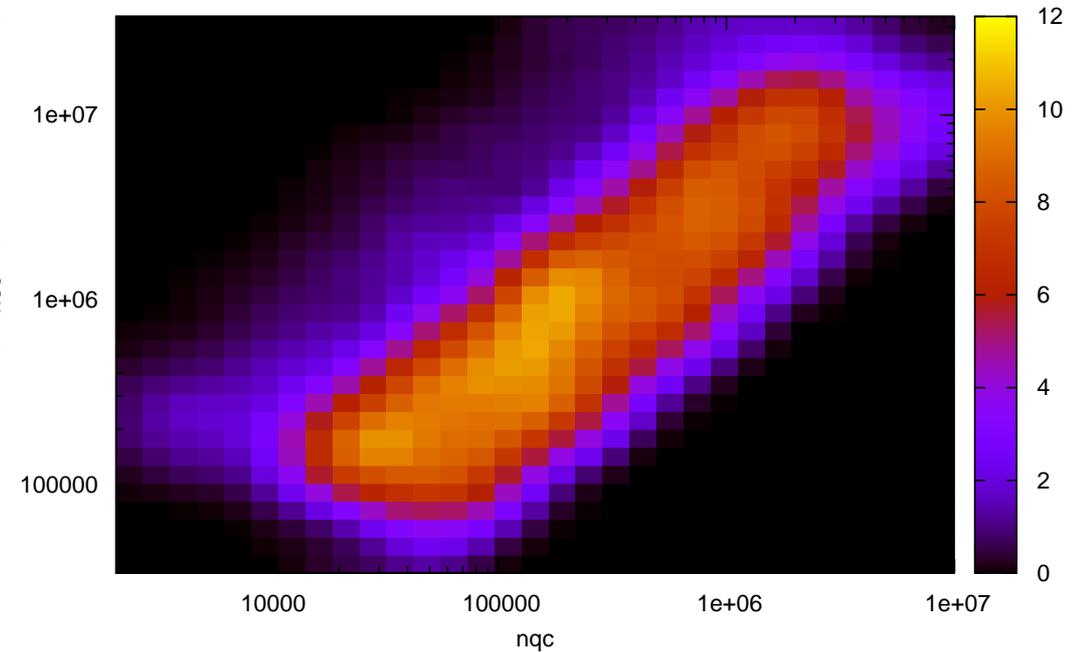


Plots: 3d: $Q : E$: histogram

raw

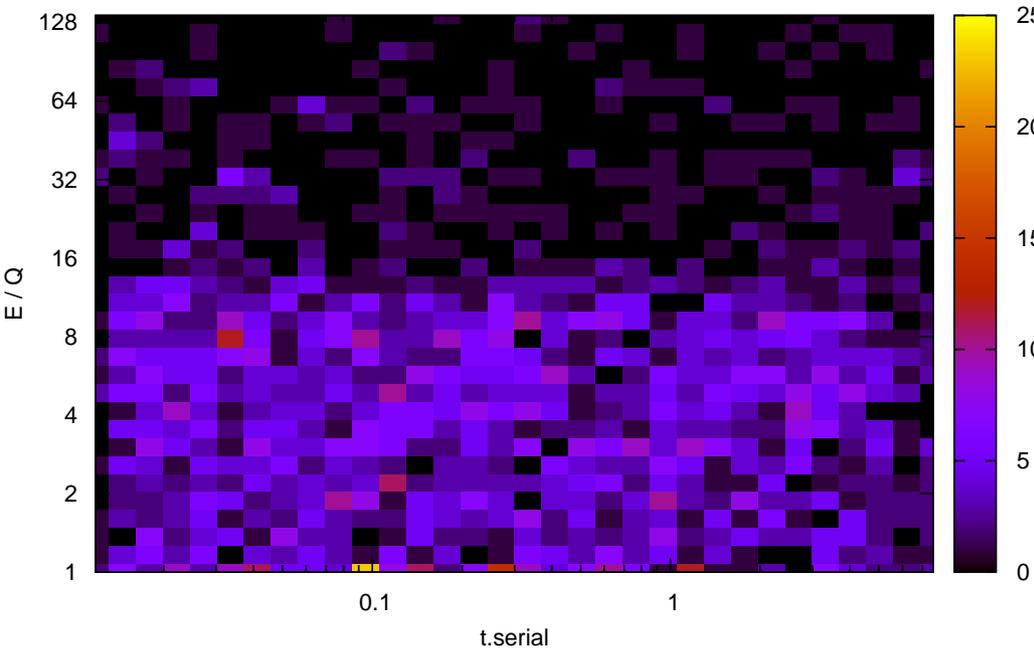


smoothed



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

raw



smoothed

