

Malleable task-graph scheduling with a practical speed-up model

Loris Marchal¹ Bertrand Simon¹ Oliver Sinnen² Frédéric Vivien¹
 1: LIP, CNRS, INRIA, ENS de Lyon and Univ. de Lyon. 2: University of Auckland, NZ.

Contact: bertrand.simon@ens-lyon.fr

Objectives

Optimize the time performance of multifrontal sparse direct solvers (e.g., MUMPS).

- ▶ Computations described by a **tree of tasks**
- ▶ Generalization to **Series-Parallel graphs** – i.e., $G = T \mid G_1; G_2 \mid G_1 \parallel G_2$

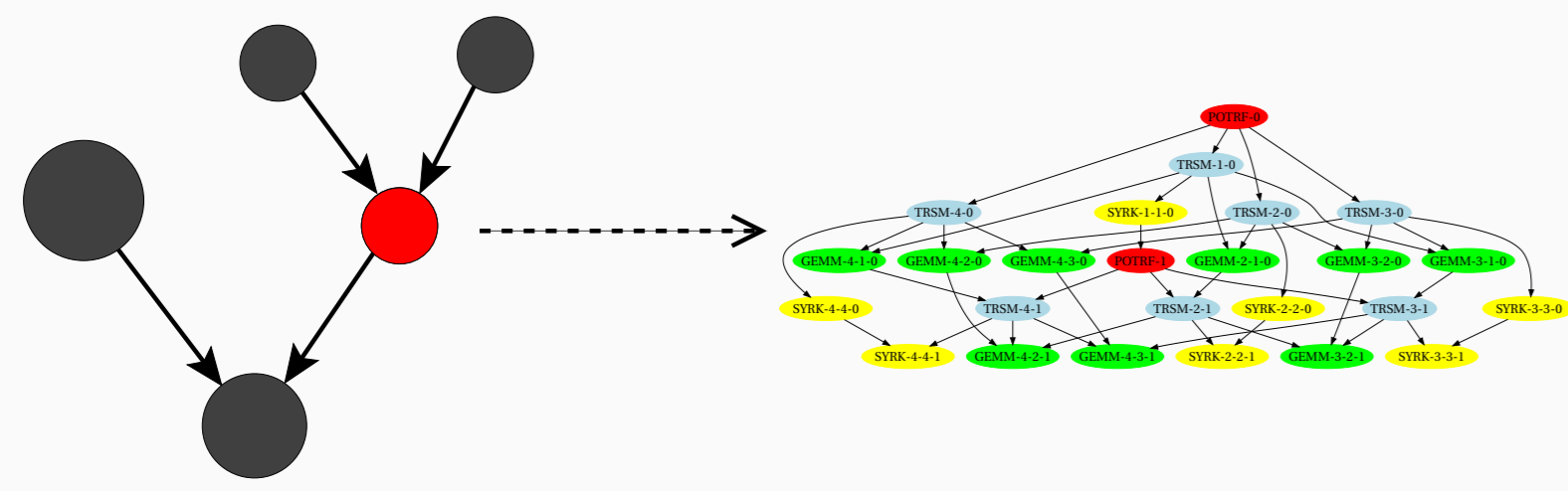
We aim at:

- ▶ Guaranteeing widely used algorithms
- ▶ Designing better scheduling algorithms

Why malleable task trees suffice?

Coarse-grain picture

- ▶ Each task: partial factorization, graph of smaller sub-tasks



Reason of this abstraction

- ▶ Expand all tasks: high complexity to schedule
- ▶ Scheduling trees simpler than general graphs

Behavior of coarse-grain tasks

- ▶ **Parallel** and **malleable**
- ▶ Speed-up model → trade-off between:
 - **Accuracy**: fits well the data
 - **Tractability**: guaranteed algorithms

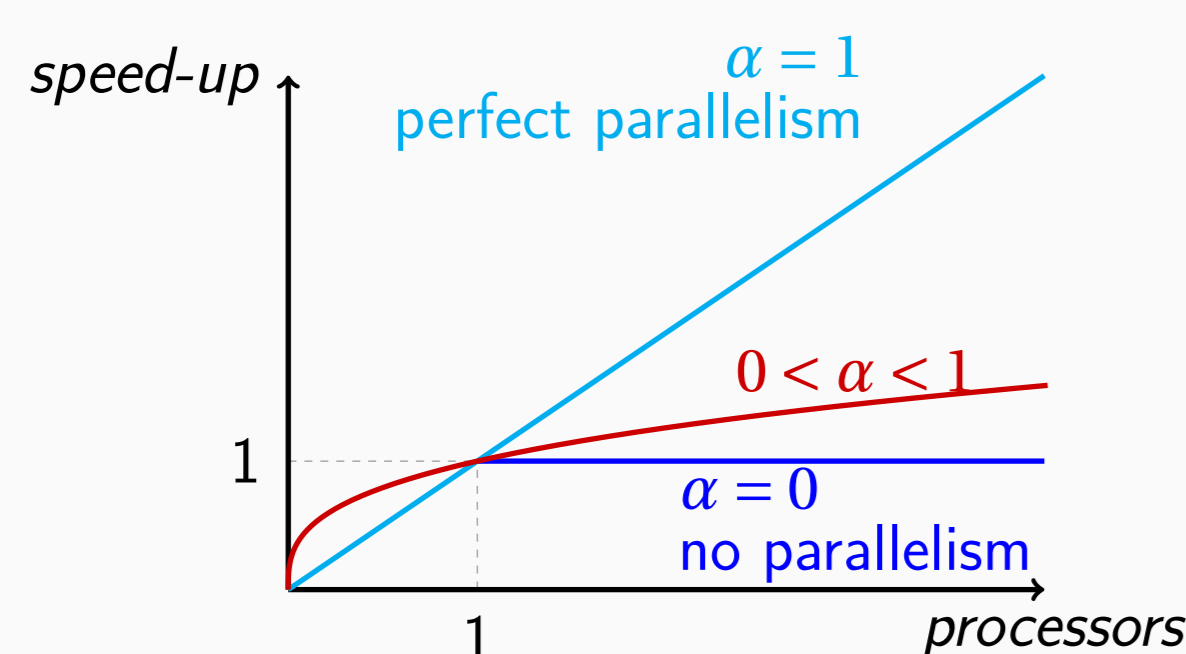
Previous work: Prasanna & Musicus model

Focus on two quantities

$$\text{speed-up}(p) = \frac{\text{time}(1 \text{ proc.})}{\text{time}(p \text{ proc.})} \mid \text{work}(p) = p \cdot \text{time}(p \text{ proc.})$$

Study model: $\text{speed-up}(p) = p^\alpha$

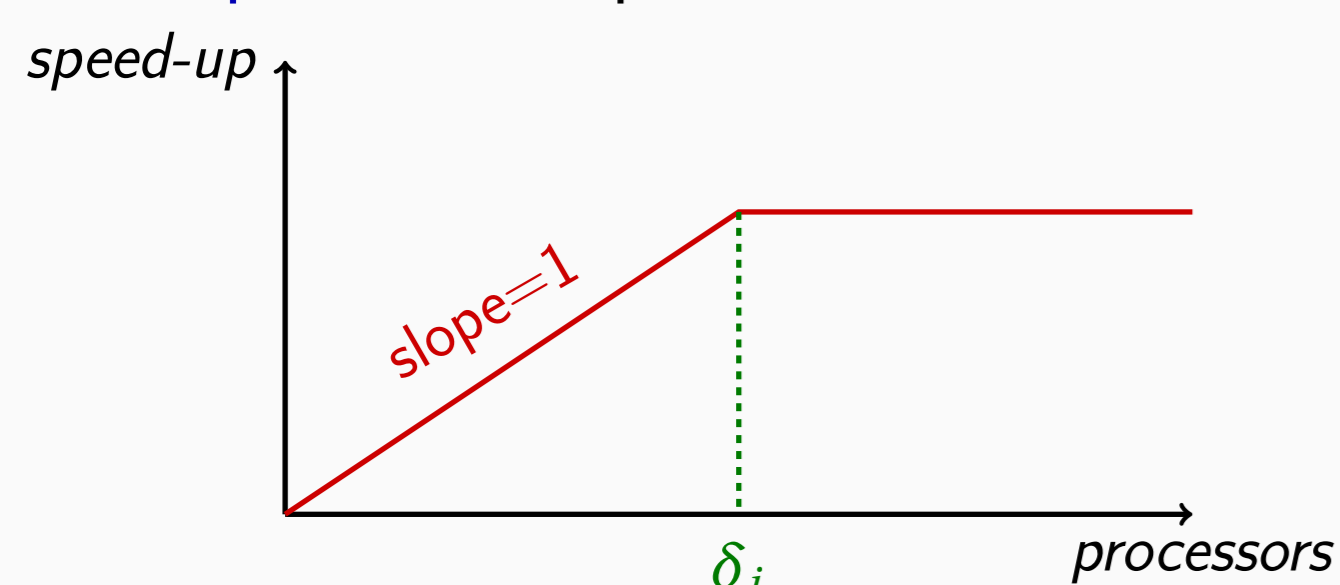
- ▶ Average Accuracy ☺
- ▶ Rational number of processors ☺
- ▶ Optimal algorithms for Series-Parallel graph ☺
- ▶ No guarantees for distributed platforms ☹
- ▶ Task finish times complex to compute ☹



A simple yet practical model

Parallel malleable tasks

- ▶ **Perfect parallelism** up to a **threshold**:



- ▶ Total work: w_i — Threshold: δ_i
- ▶ Rational allocation for free (McNaughton's wrap-around rule)

Related work

Non-increasing speed-up and work

- ▶ Independent tasks: theoretical FPTAS and practical 2-approximations [Jansen 2004, Fan et al. 2012]
- ▶ SP-graphs: ≈ 2.6 -approximation [Lepère et al. 2001]. With concave speed-up: $(2 + \epsilon)$ -approximation of unspecified complexity [Makarychev et al. 2014]

Specific speed-up function

- ▶ Same model: 2-approximation [Balmin et al. 2013] named FLOWFLEX (see experimental setup)
- ▶ [Kell et al. 2015]: $\text{time} = \frac{w_i}{p} + (p-1)c$; 2-approximation for $p=3$, open for $p \geq 4$

NP-Completeness of the problem

Complexity depending on the model

- ▶ Malleability + perfect parallelism $\Rightarrow P$
- ▶ Adding thresholds $\Rightarrow NP$ -complete
- ▶ Arguably complex proof [Drozdowski and Kubiak 1999]

Contribution

- ▶ **New NP-completeness proof**

The widely used PROPMAPPING

Simple allocation for trees or SP-graphs

- ▶ On a series composition $G = (G_1; G_2)$: give all available processors to G_1 , then to G_2
- ▶ On $(G_1 \parallel G_2)$: give a **constant share** to G_i , **proportional to its weight w_i**
- ▶ Algorithm on graph G with q processors:

PROPMAPPING (G, q)

if $G = G_1; G_2; \dots; G_k$ then $\forall i, p_i \leftarrow q$ if $G = G_1 \parallel G_2 \parallel \dots \parallel G_k$ then $\forall i, p_i = q w_i / \sum w_j$

Call PROPMAPPING (G_i, p_i) for each G_i

- ▶ Then schedule each task on p_i processors as soon as it is ready

Notes

- ▶ Moldable schedule (constant allocation)
- ▶ Unaware of task thresholds

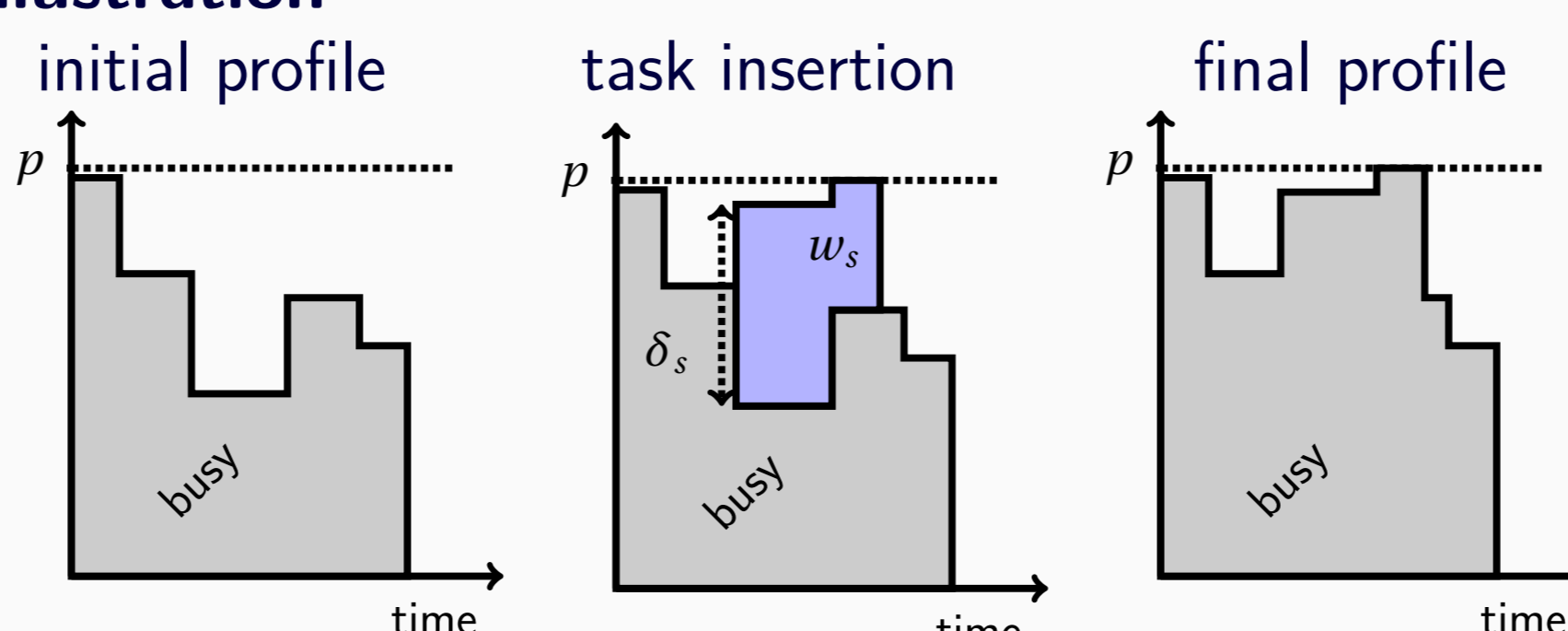
Theorem: PROPMAPPING is a 2-approximation.

A new strategy: GREEDY-FILLING

Algorithm

- ▶ Assign priorities to tasks (usually bottom-level)
- ▶ Consider free tasks by decreasing priority
- ▶ **Greedy insert** each task in the schedule:
 - Compute the earliest starting time
 - **Pour** task into the **available processor space**, respecting thresholds

Illustration



Theorem: GREEDY-FILLING is a 2-approximation.

Experimental setup

Third algorithm for comparison: FLOWFLEX

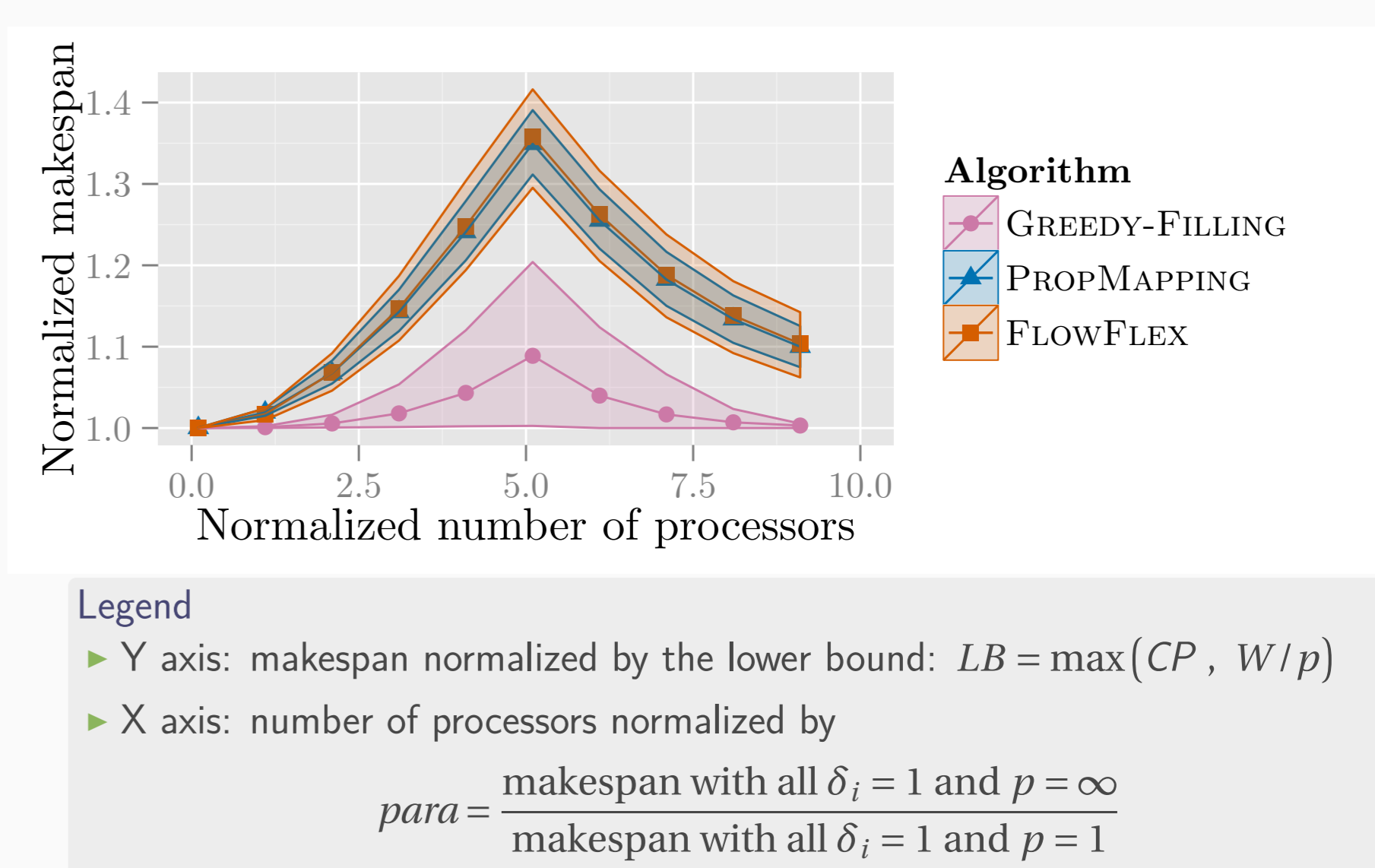
- ▶ 2-approximation from [Balmin et al. 2013] to schedule *Malleable Flows of MapReduce Jobs*
- ▶ Solve the problem on an infinite number of processors
- ▶ **Downscale the allocation** on intervals when it is needed

Two datasets

- ▶ SYNTH: synthetic SP-graphs with $\delta_i = \alpha \times w_i$
- ▶ TREES: assembly trees of sparse matrices, $\delta_i = \alpha \times w_i$.

Validation of GREEDY-FILLING

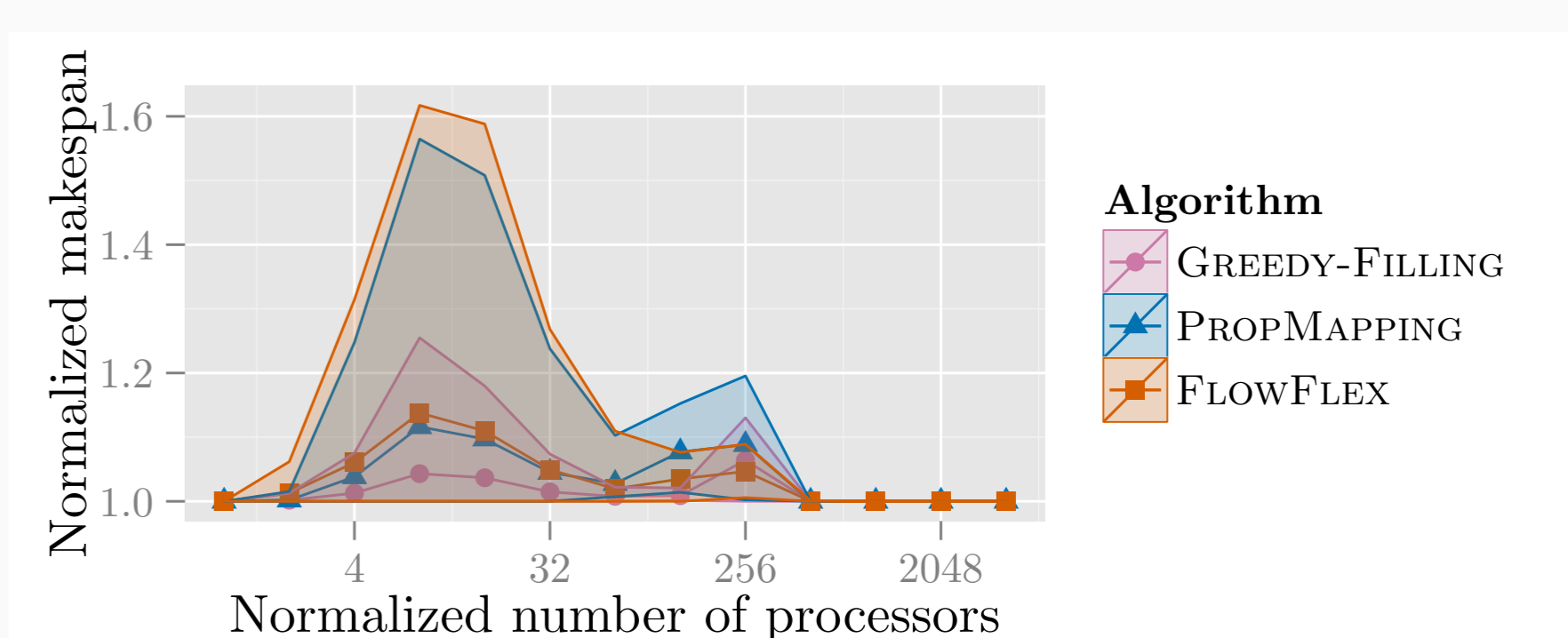
Results on SYNTH



Comments

- ▶ Plot: mean + ribbon with 90% of the results
- ▶ Small/large number of processors: similar results as the problem is simple
- ▶ **GREEDY-FILLING**: $\approx 25\%$ of gain < 20% from the LB

Results on TREES



Comments

- ▶ Results shape depends a lot on the matrix
- ▶ Here: one matrix with different ordering and amalgamation parameters
- ▶ GREEDY-FILLING is (almost always) better than both others
- ▶ Smaller maximum gain (around 15%)

Conclusion

On the algorithms

- ▶ PROPMAPPING: does not take advantage of malleability
- ▶ FLOWFLEX: produces gaps that cannot be filled afterwards
- ▶ GREEDY-FILLING: simple, greedy, close to the lower bound

On the model

- ▶ Simplest model to account for limited parallelism
- ▶ Still NP-complete
- ▶ Possible to derive theoretical guarantees (2-approximation algorithms)