

6th SOS Workshop on Distributed Supercomputing: Data Intensive Computing
March 4-6, 2002, Badhotel Bristol, Leukerbad, Valais, Switzerland

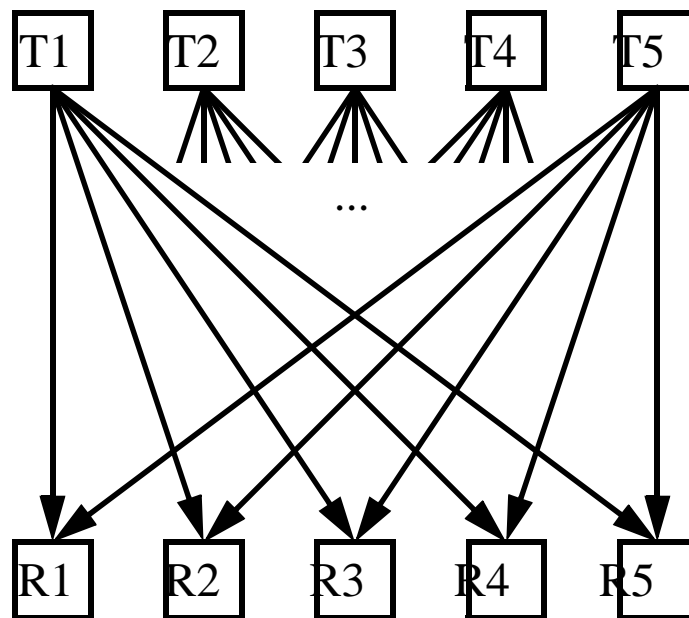
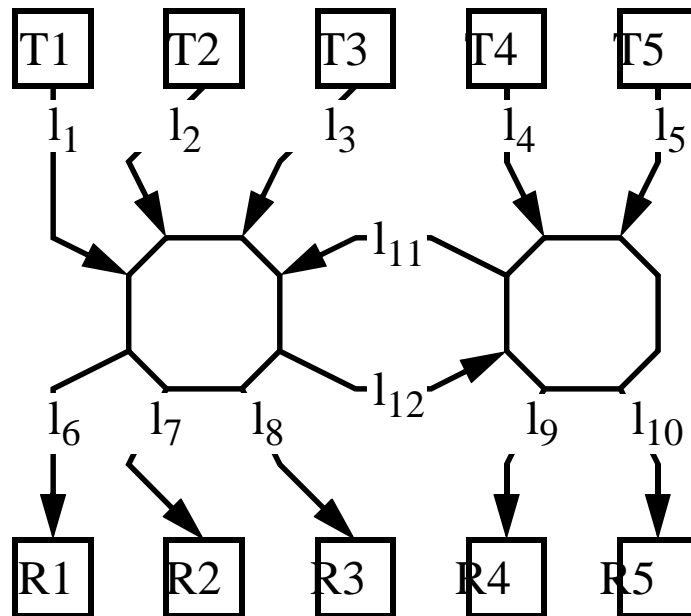
Network Topology-aware Traffic Scheduling

Emin Gabrielyan

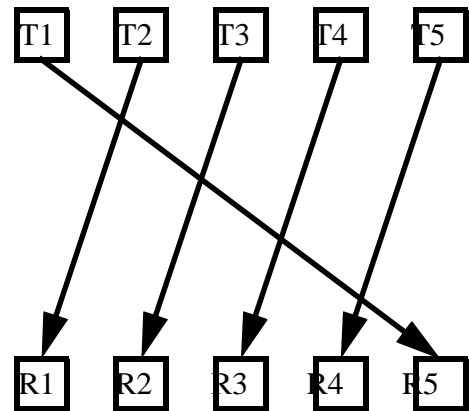
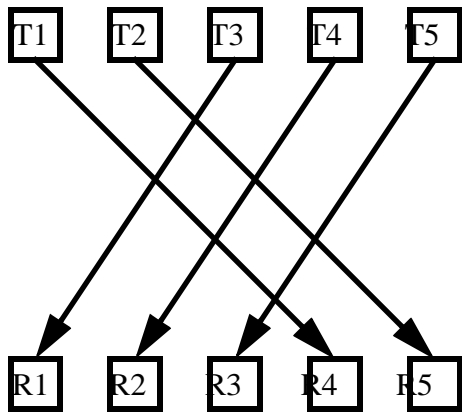
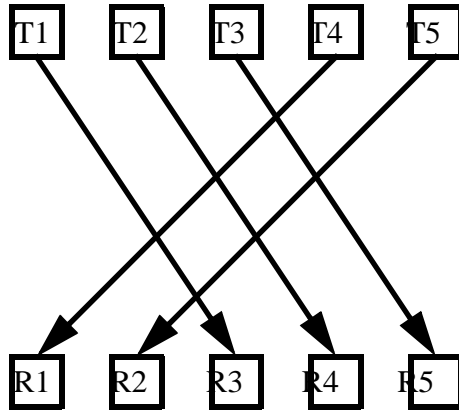
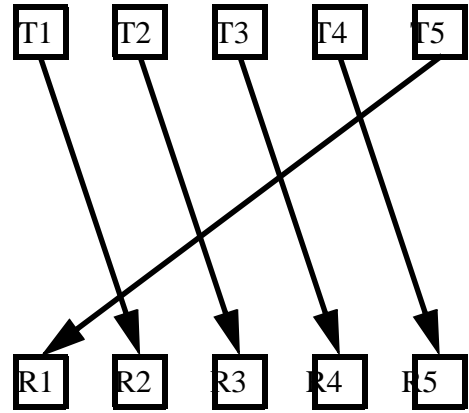
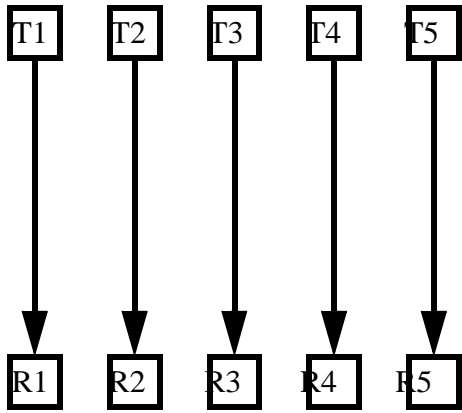
*École Polytechnique Fédérale de
Lausanne, Switzerland*

Emin.Gabrielyan@epfl.ch

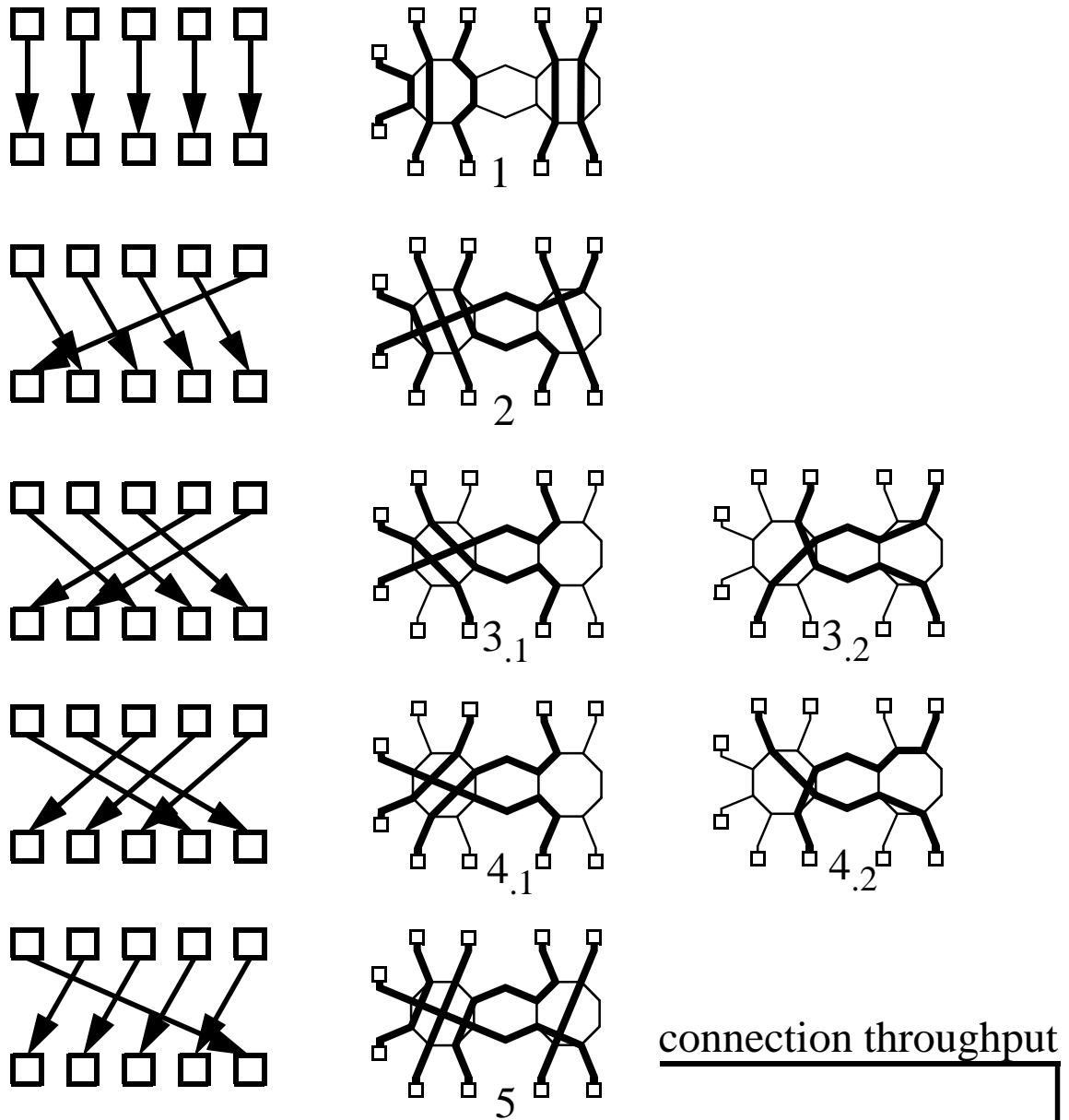
25-transfer data exchange



Round-robin schedule



Round-robin Throughput



mean number of connections per timeframe

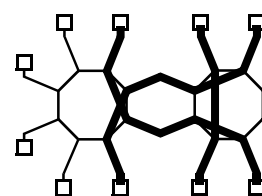
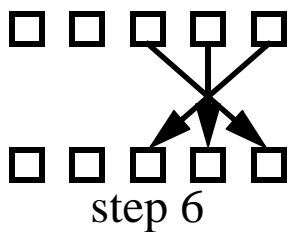
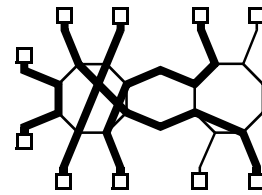
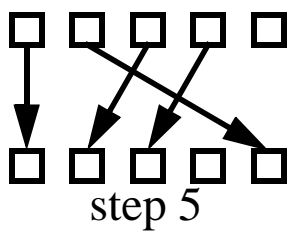
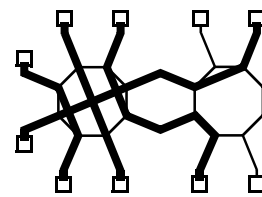
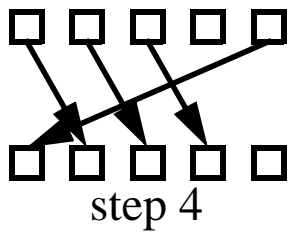
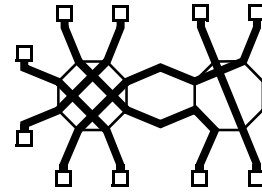
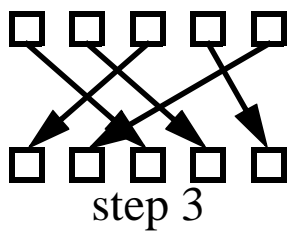
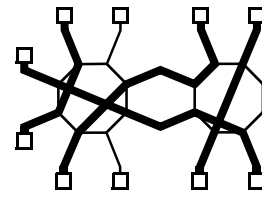
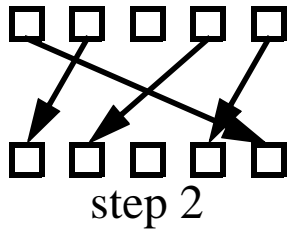
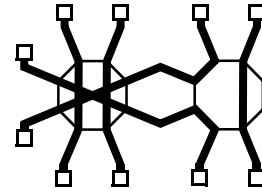
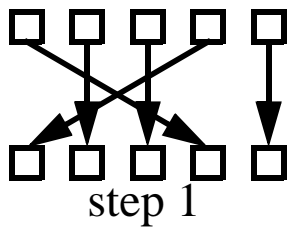
$$T_{\text{roundrobin}} = \frac{25}{7} \cdot 100\text{MB/s} = 357\text{MB/s}$$

total throughput ←

number of transfers ←

number of timeframes ←

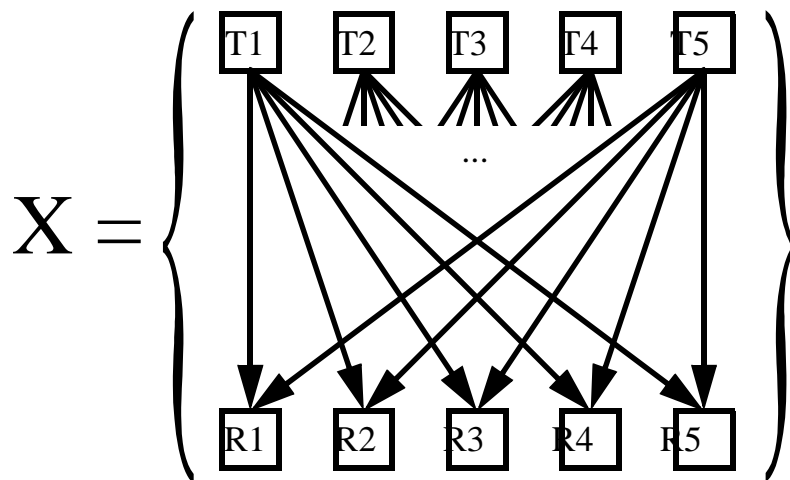
Liquid Schedule



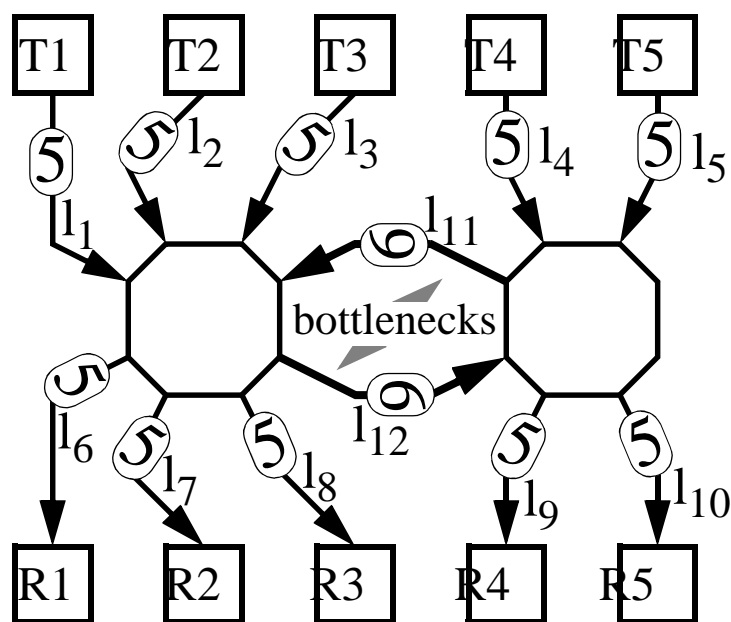
$$T_{liquid} = \underbrace{25/6} \cdot 100MB/s = 416MB/s$$

mean number of connections per step

Load of Links and Transfers



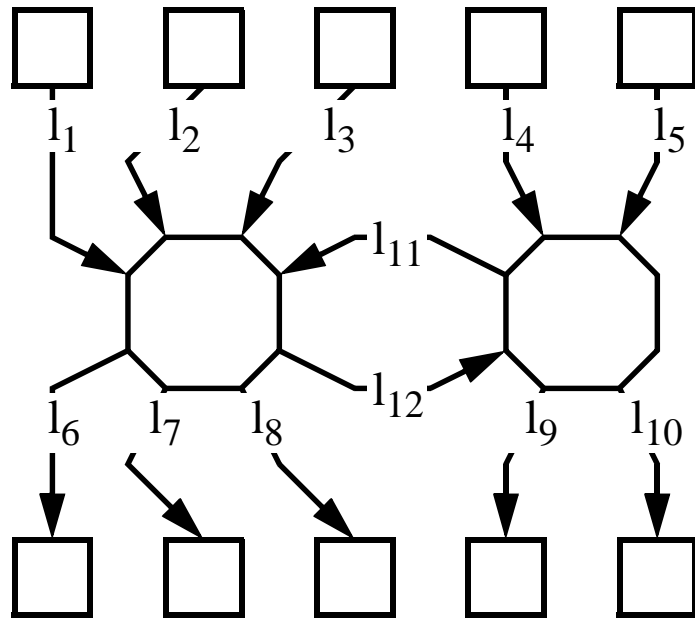
The 25 transfer traffic



$$\lambda(l_1, X) = 5, \dots, \lambda(l_{12}, X) = 6$$

Transfers: $\{l_1, l_6\}, \dots, \{l_1, l_{12}, l_6\}, \dots$

Duration of the Traffic



$$X = \left\{ \begin{array}{l} \{l_1, l_6\}, \{l_1, l_7\}, \{l_1, l_8\}, \{l_1, l_{12}, l_9\}, \{l_1, l_{12}, l_{10}\}, \\ \{l_2, l_6\}, \{l_2, l_7\}, \{l_2, l_8\}, \{l_2, l_{12}, l_9\}, \{l_2, l_{12}, l_{10}\}, \\ \{l_3, l_6\}, \{l_3, l_7\}, \{l_3, l_8\}, \{l_3, l_{12}, l_9\}, \{l_3, l_{12}, l_{10}\}, \\ \{l_4, l_{11}, l_6\}, \{l_4, l_{11}, l_7\}, \{l_4, l_{11}, l_8\}, \{l_4, l_9\}, \{l_4, l_{10}\}, \\ \{l_5, l_{11}, l_6\}, \{l_5, l_{11}, l_7\}, \{l_5, l_{11}, l_8\}, \{l_5, l_9\}, \{l_5, l_{10}\} \end{array} \right\}$$

$$\lambda(l_1, X) = 5, \lambda(l_2, X) = 5, \dots$$

$$\lambda(l_{11}, X) = 6, \lambda(l_{12}, X) = 6$$

$$\Lambda(X) = 6$$

Liquid Throughput

$$X = \left\{ \begin{array}{l} \{l_1, l_6\}, \{l_1, l_7\}, \{l_1, l_8\}, \{l_1, l_{12}, l_9\}, \{l_1, l_{12}, l_{10}\}, \\ \{l_2, l_6\}, \{l_2, l_7\}, \{l_2, l_8\}, \{l_2, l_{12}, l_9\}, \{l_2, l_{12}, l_{10}\}, \\ \{l_3, l_6\}, \{l_3, l_7\}, \{l_3, l_8\}, \{l_3, l_{12}, l_9\}, \{l_3, l_{12}, l_{10}\}, \\ \{l_4, l_{11}, l_6\}, \{l_4, l_{11}, l_7\}, \{l_4, l_{11}, l_8\}, \{l_4, l_9\}, \{l_4, l_{10}\}, \\ \{l_5, l_{11}, l_6\}, \{l_5, l_{11}, l_7\}, \{l_5, l_{11}, l_8\}, \{l_5, l_9\}, \{l_5, l_{10}\} \end{array} \right\}$$

the throughput of a single link

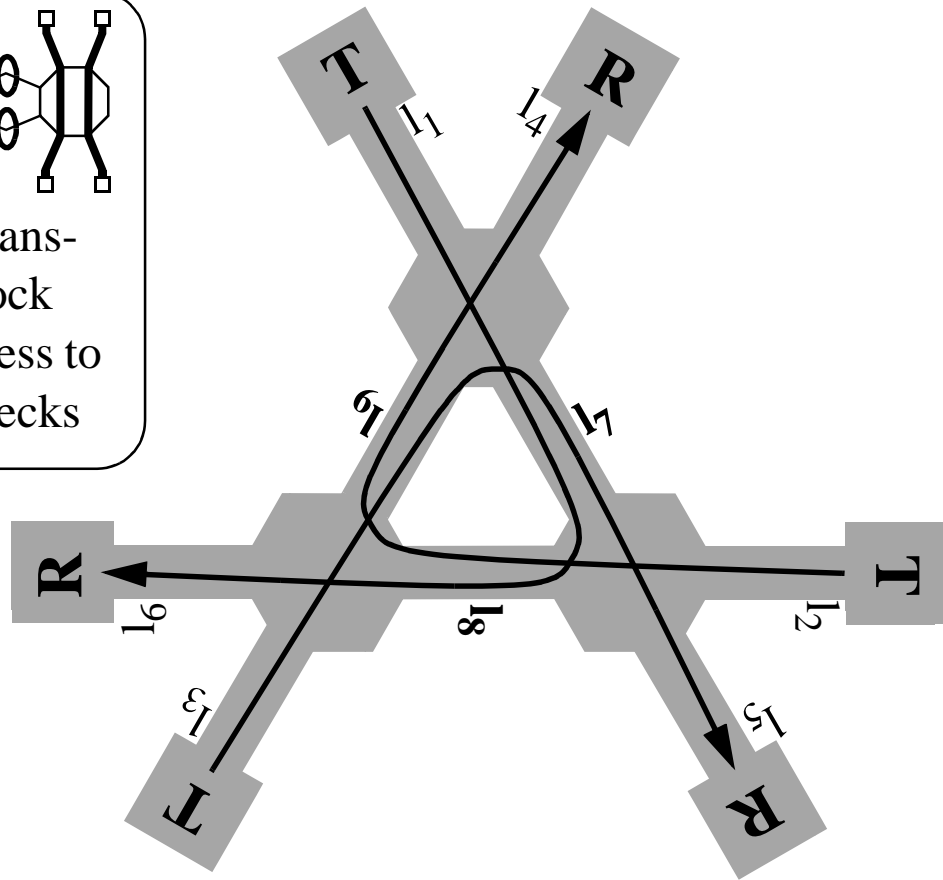
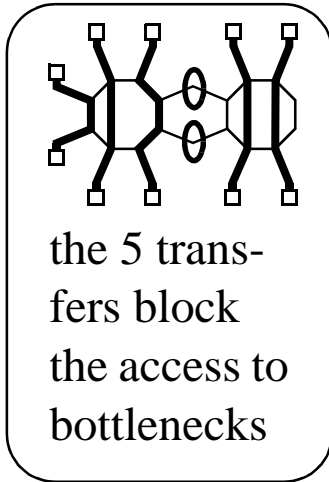
total number of transfers

$$T_{liquid} = \frac{\#(X)}{\Lambda(X)} \cdot T_{link} =$$

the duration of the traffic (the load of its bottlenecks)

$$= \frac{25}{6} \cdot 100MB/s = 417MB/s$$

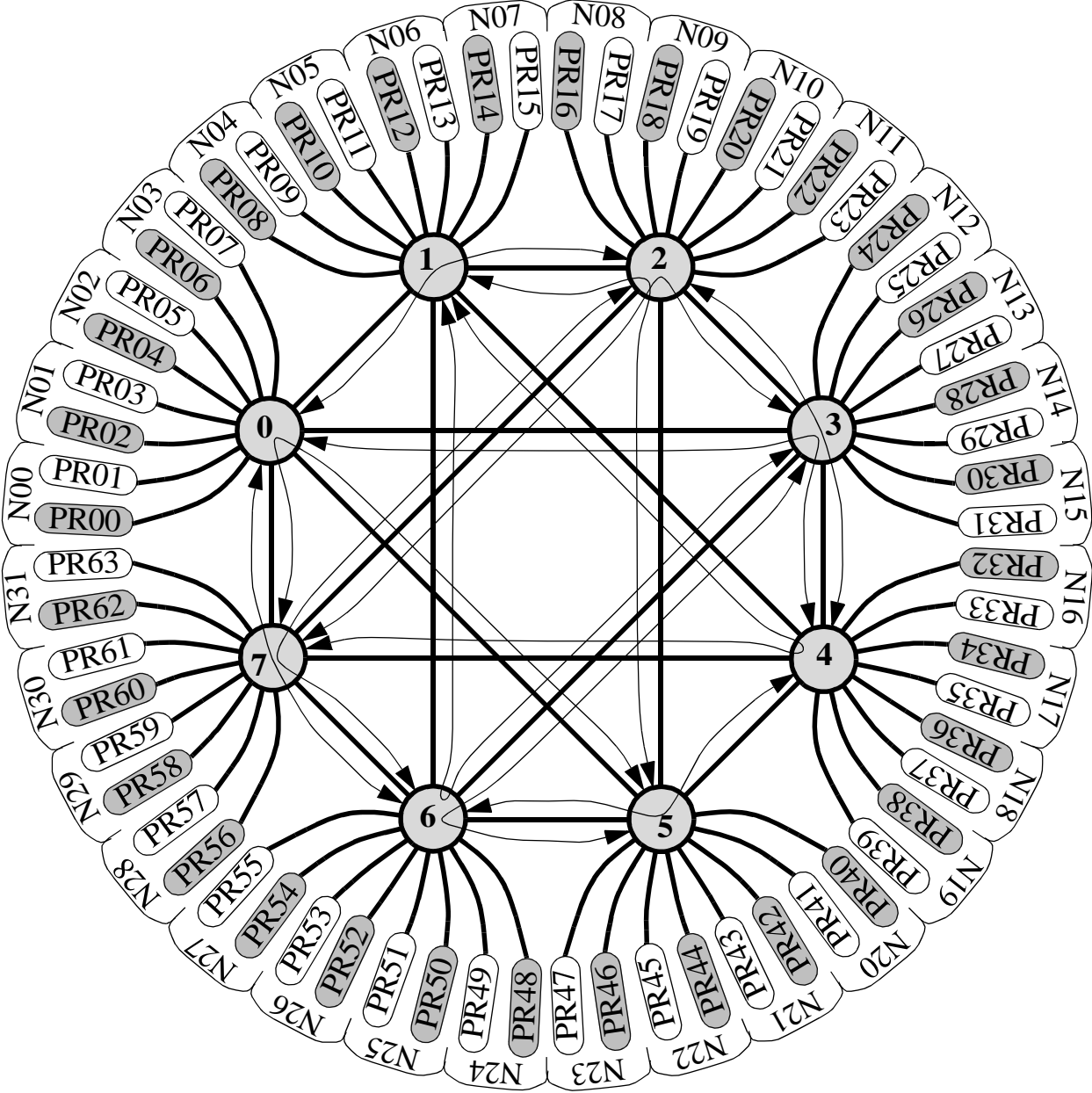
No liquid schedule



$$X = \left\{ \begin{array}{l} \{l_1, l_7, l_8, l_6\}, \\ \{l_2, l_8, l_9, l_4\}, \\ \{l_3, l_9, l_7, l_5\} \end{array} \right\} \quad \begin{array}{l} \#(X) = 3 \\ \Lambda(X) = 2 \end{array}$$

$$\begin{aligned} T_{liquid} &= \frac{\#(X)}{\Lambda(X)} \cdot T_{link} = \\ &= 3/2 \cdot 100MB/s = 150MB/s \end{aligned}$$

Swiss-T1 Cluster



PR01 Receiving Processor

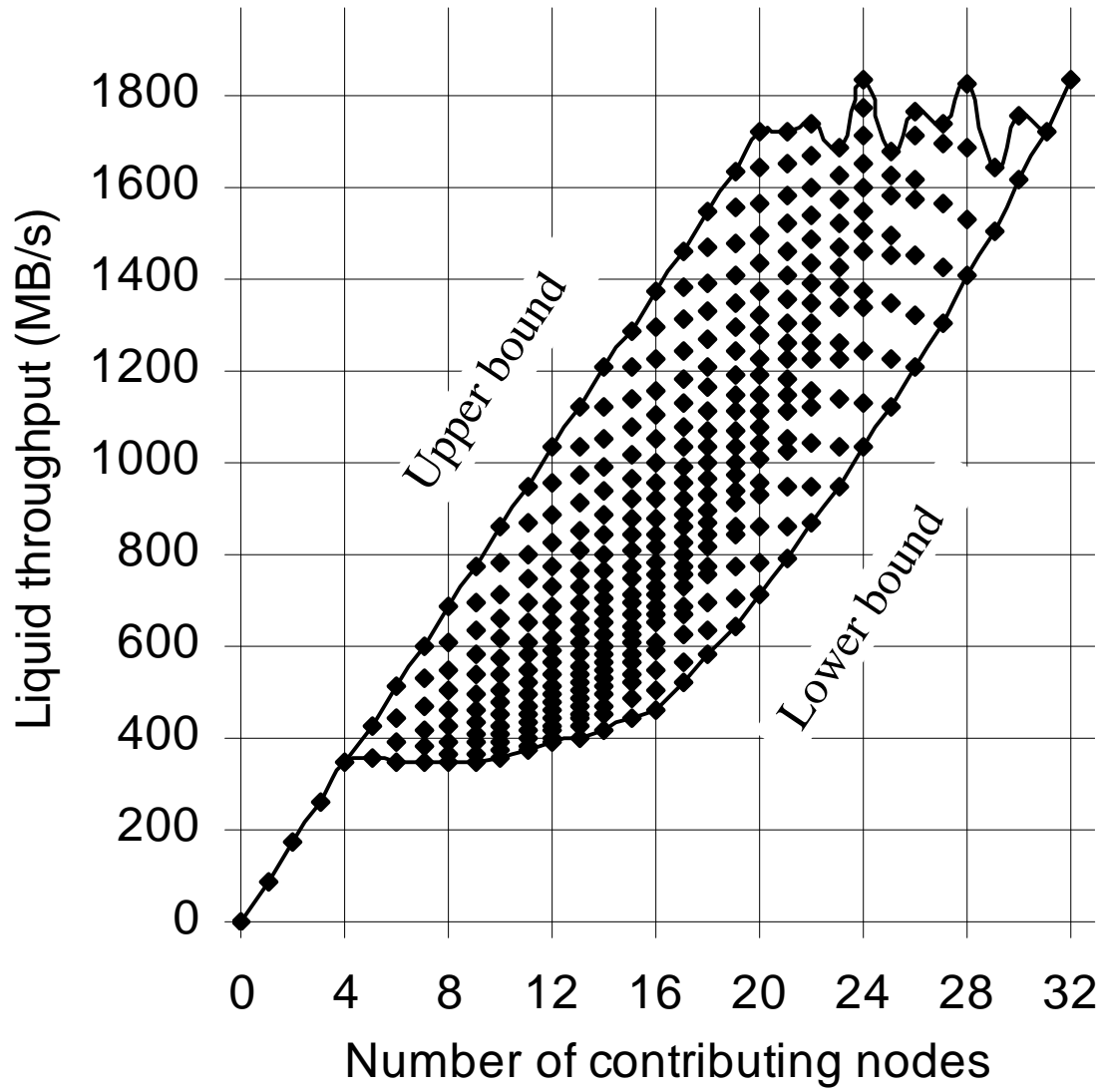
PR00 Sending Processor

↑ Routing information
 ↓ Network link

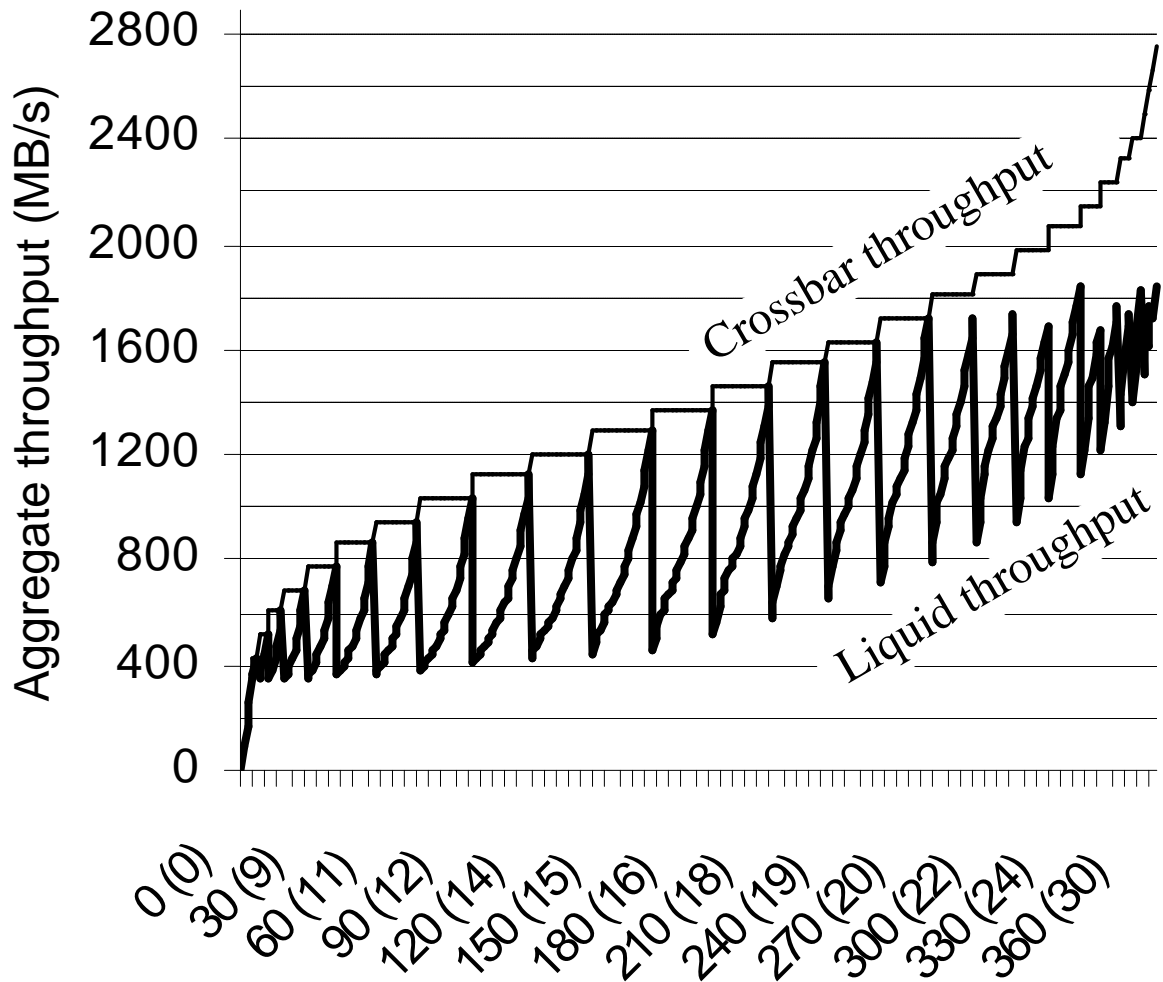
N00 Node

0 Switch

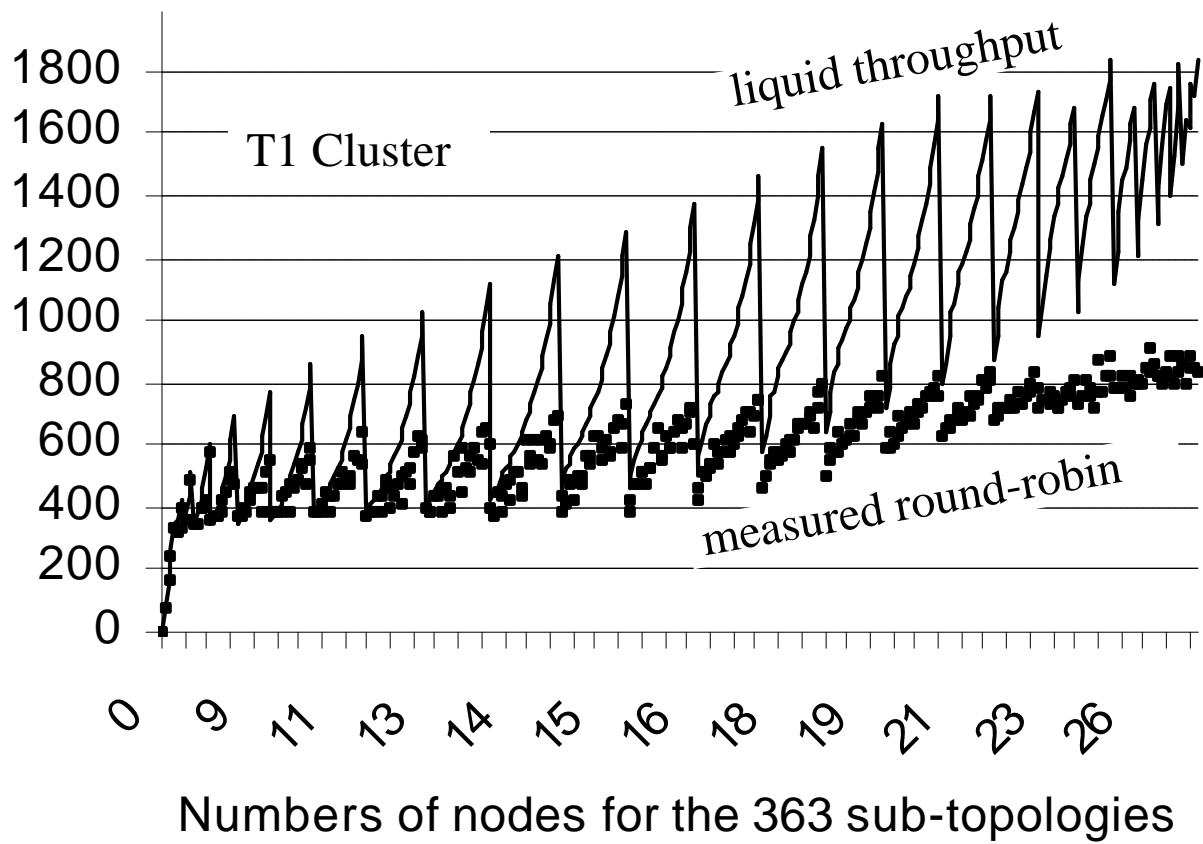
363 Test Traffics



363-Topology Test-bed



Round-robin throughput



Team: set of non-congesting transfers
using all bottlenecks

$$X = \left\{ \begin{array}{l} \{1_1, 1_6\}, \{1_1, 1_7\}, \{1_1, 1_8\}, \{1_1, \mathbf{1}_{12}, 1_9\}, \{1_1, \mathbf{1}_{12}, 1_{10}\}, \\ \{1_2, 1_6\}, \{1_2, 1_7\}, \{1_2, 1_8\}, \{1_2, \mathbf{1}_{12}, 1_9\}, \{1_2, \mathbf{1}_{12}, 1_{10}\}, \\ \{1_3, 1_6\}, \{1_3, 1_7\}, \{1_3, 1_8\}, \{1_3, \mathbf{1}_{12}, 1_9\}, \{1_3, \mathbf{1}_{12}, 1_{10}\}, \\ \{1_4, \mathbf{1}_{11}, 1_6\}, \{1_4, \mathbf{1}_{11}, 1_7\}, \{1_4, 1_{11}, 1_8\}, \{1_4, 1_9\}, \{1_4, 1_{10}\}, \\ \{1_5, \mathbf{1}_{11}, 1_6\}, \{1_5, \mathbf{1}_{11}, 1_7\}, \{1_5, 1_{11}, 1_8\}, \{1_5, 1_9\}, \{1_5, 1_{10}\} \end{array} \right\}$$

$$\alpha = \left\{ \begin{array}{l} \left(\begin{array}{l} \{1_1, \mathbf{1}_{12}, 1_9\}, \\ \{1_2, 1_7\}, \\ \{1_3, 1_8\}, \\ \{1_4, \mathbf{1}_{11}, 1_6\}, \\ \{1_5, 1_{10}\} \end{array} \right), \left(\begin{array}{l} \{1_1, \mathbf{1}_{12}, 1_{10}\}, \\ \{1_2, 1_6\}, \\ \{1_4, \mathbf{1}_{11}, 1_7\}, \\ \{1_5, 1_9\} \end{array} \right), \left(\begin{array}{l} \{1_1, 1_8\}, \\ \{1_2, \mathbf{1}_{12}, 1_9\}, \\ \{1_3, 1_6\}, \\ \{1_4, 1_{10}\}, \\ \{1_5, \mathbf{1}_{11}, 1_7\} \end{array} \right), \\ \left(\begin{array}{l} \{1_1, 1_7\}, \\ \{1_2, 1_8\}, \\ \{1_3, \mathbf{1}_{12}, 1_9\}, \\ \{1_5, \mathbf{1}_{11}, 1_6\} \end{array} \right), \left(\begin{array}{l} \{1_1, 1_6\}, \\ \{1_2, \mathbf{1}_{12}, 1_{10}\}, \\ \{1_3, 1_7\}, \\ \{1_4, \mathbf{1}_{11}, 1_8\} \end{array} \right), \left(\begin{array}{l} \{1_3, \mathbf{1}_{12}, 1_{10}\}, \\ \{1_4, 1_9\}, \\ \{1_5, 1_{11}, 1_8\} \end{array} \right) \end{array} \right\}$$

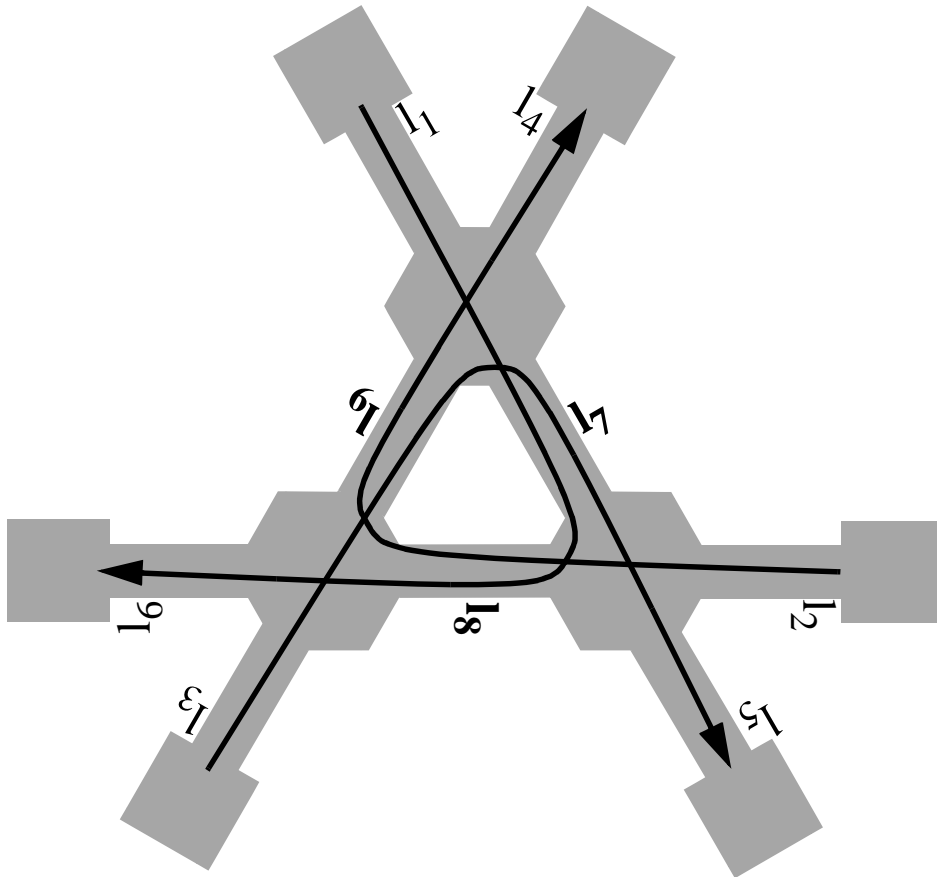
schedule α is liquid \Leftrightarrow

number of steps \swarrow \searrow load of the bottlenecks

$$\Leftrightarrow \#(\alpha) = \Lambda(X) \Leftrightarrow$$

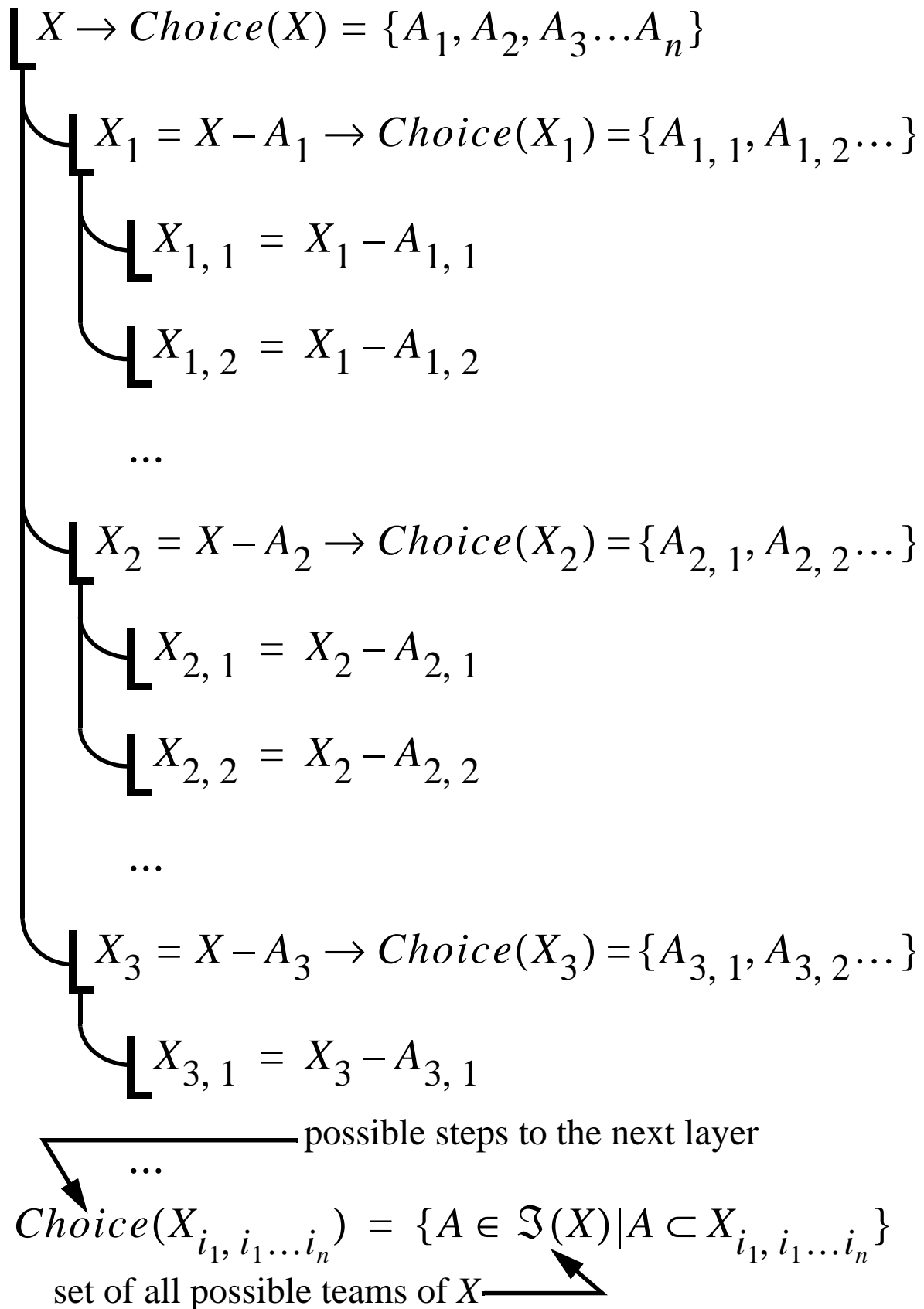
$$\Leftrightarrow \forall (A \in \alpha) A \text{ is a team of } X$$

Traffic without a team

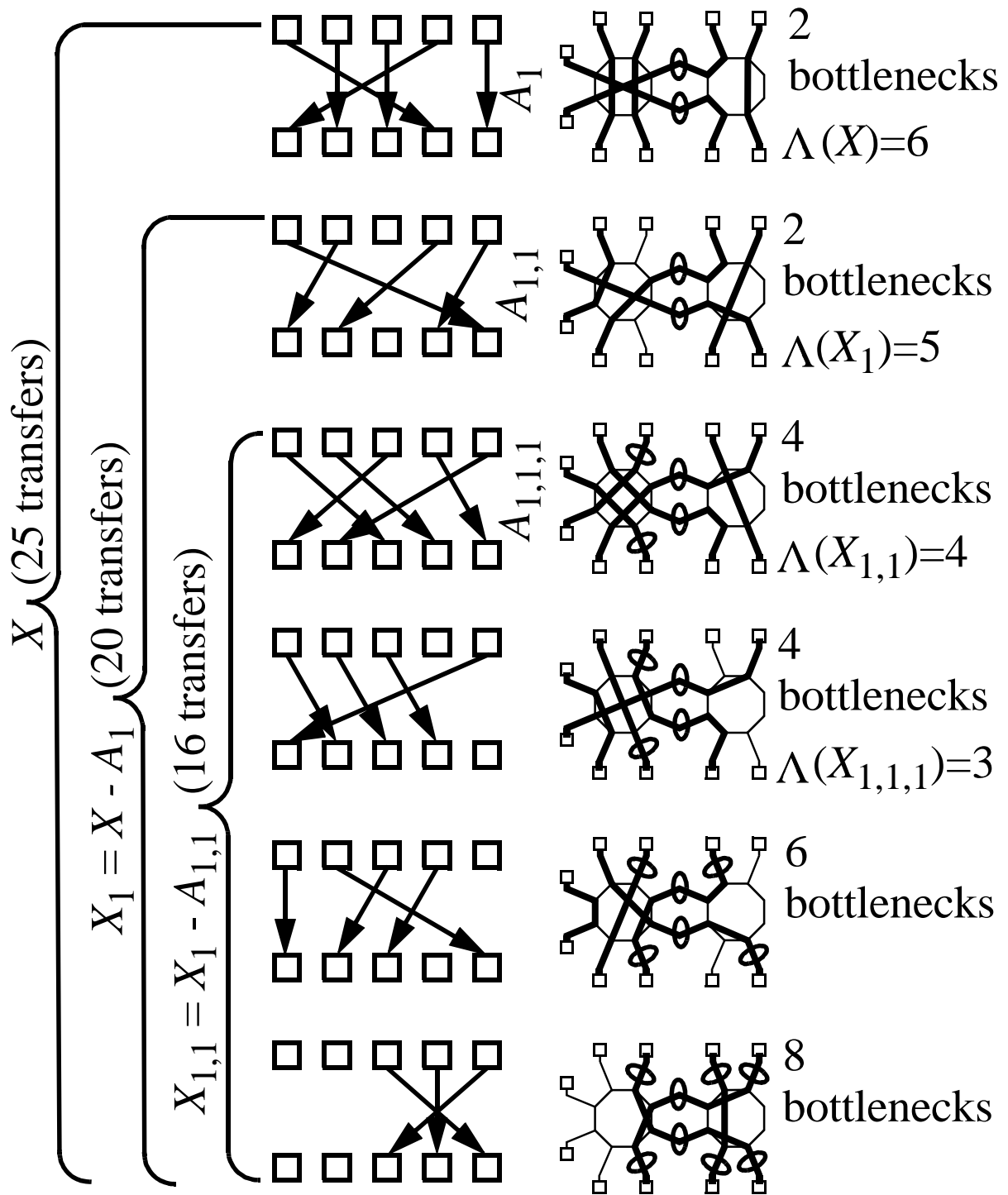


$$X = \left\{ \begin{array}{l} \{1_1, 1_7, 1_8, 1_6\}, \\ \{1_2, 1_8, 1_9, 1_4\}, \\ \{1_3, 1_9, 1_7, 1_5\} \end{array} \right\}$$

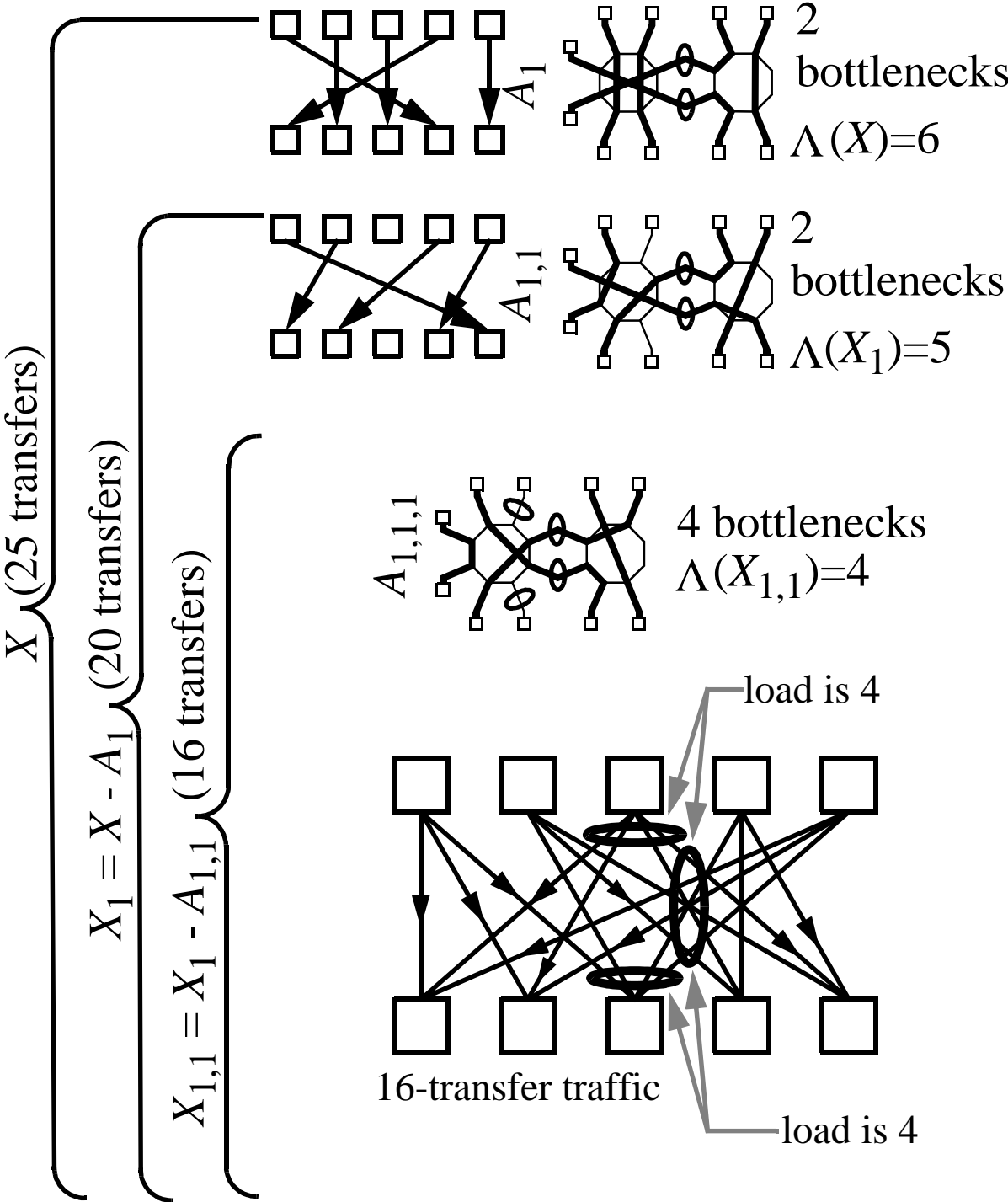
Liquid schedule search tree



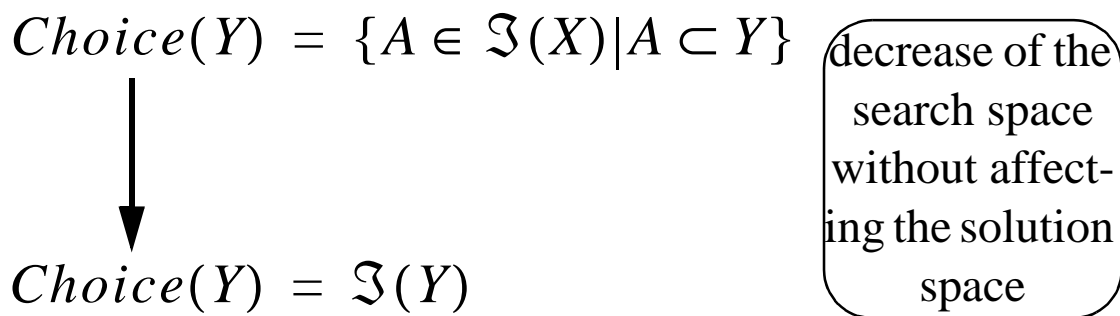
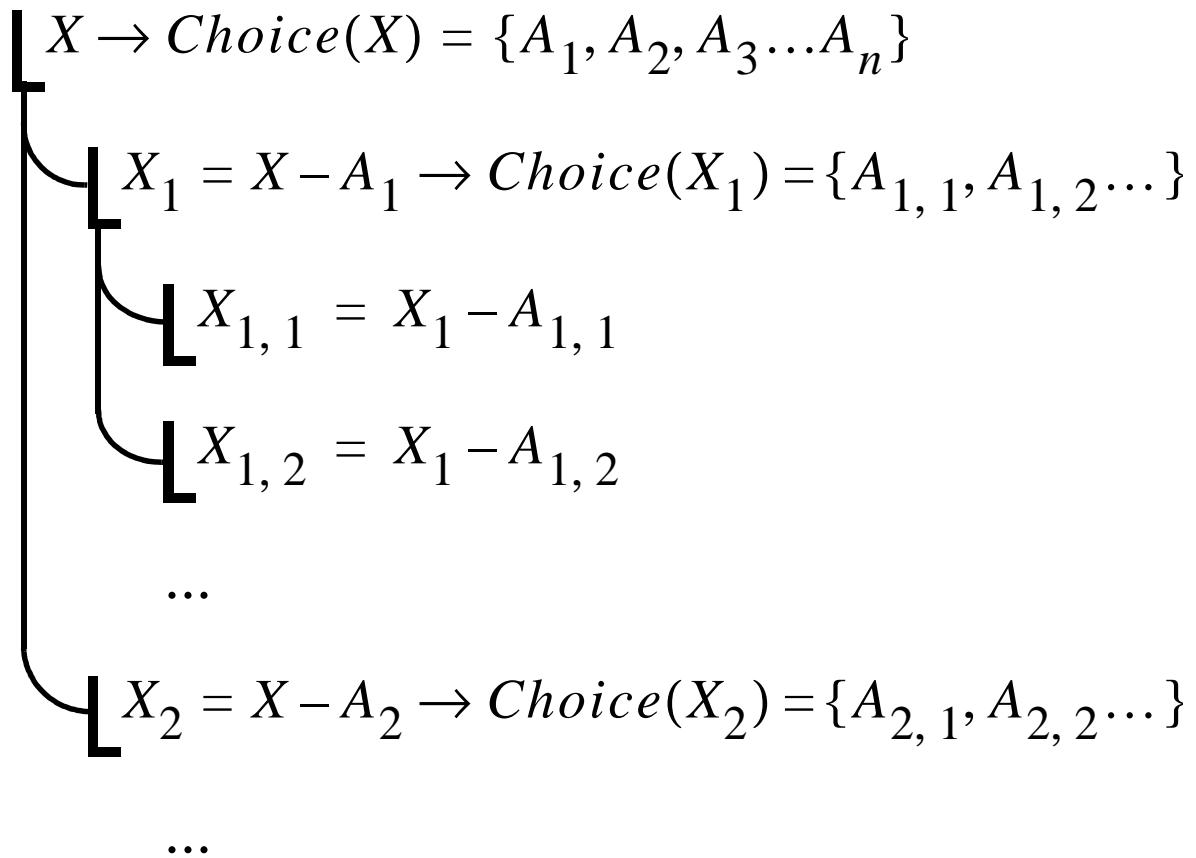
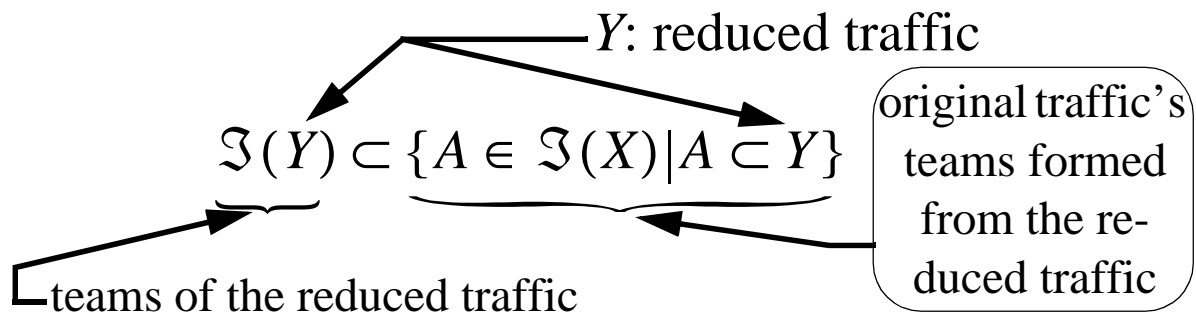
Additional bottlenecks



Prediction of Dead-ends



Liquid schedule search optimization



Liquid schedules construction

$$\underbrace{\mathfrak{S}^{full}(Y) \subset \mathfrak{S}(Y)}_{\text{full teams of the reduced traffic}}$$

$$Choice(Y) = \mathfrak{S}(Y)$$

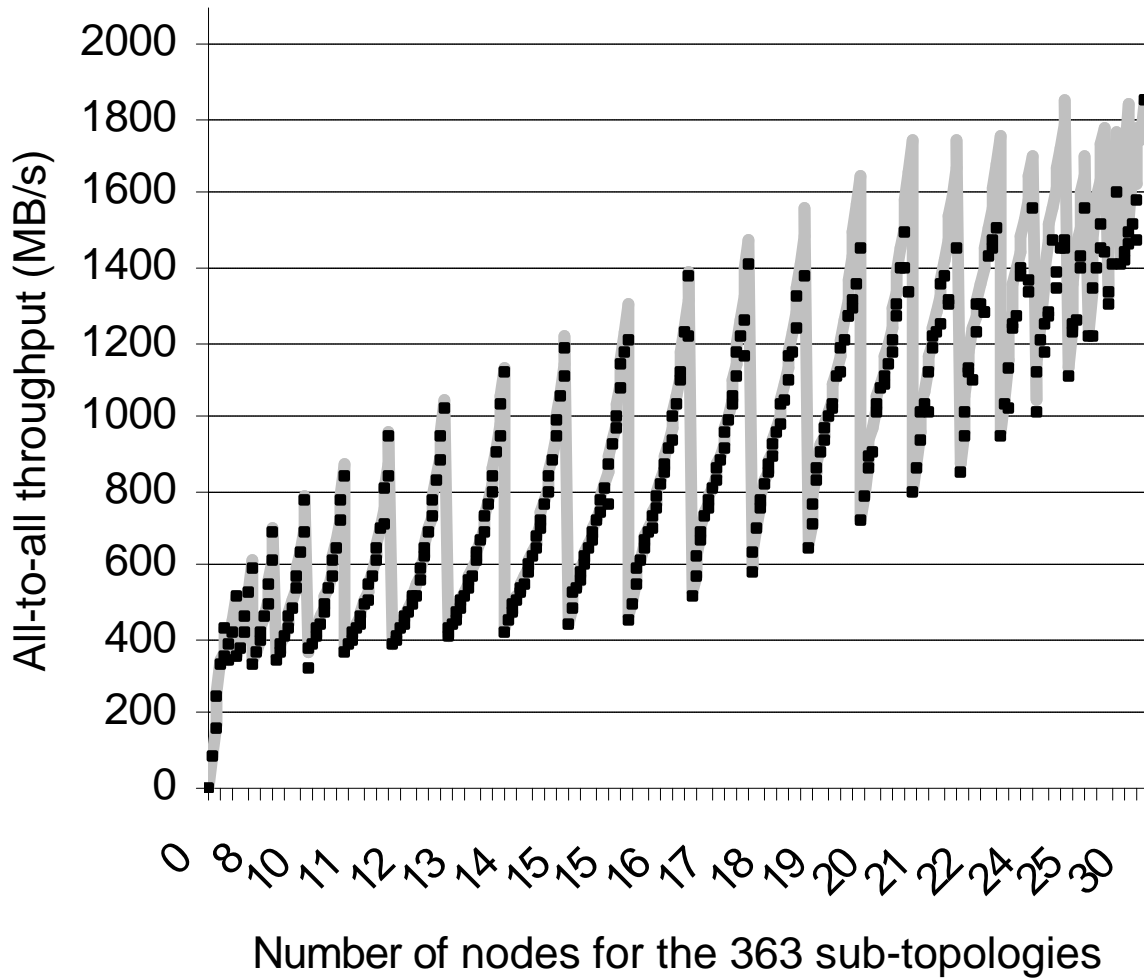


$$Choice(Y) = \mathfrak{S}^{full}(Y)$$

decrease of
the search
space without
affecting the
solution space

- For more than 90% of the test-bed topologies the search of liquid schedules took less than 0.1s on a single 500MHz processor.
- For 8 topologies out of 363 solution was not found within 24 hours.

Results



Conclusion

- Data exchanges relying on the liquid schedules may be carried out several times faster compared with topology-unaware schedules.
- Our method may be applied to applications requiring high network efficiency, such as video or voice traffic management, high energy physics data acquisition and event assembling.
- At the present we consider only static routing scheme. Dynamic routing could possibly be also combined in the algorithms.
- Fixed packet size transfers are considered.
- The network latency are neglected in comparison with the transfer times.

Thank You!

Contact: *Emin.Gabrielyan@epfl.ch*