

# Scalable and Fault Tolerant Orthogonalization Based on Randomized Aggregation

**W. N. Gansterer**, G. Niederbrucker, H. Strakova  
and S. Schulze Grotthoff

Research Lab Computational Technologies and Applications  
Research Group Theory and Applications of Algorithms  
University of Vienna, Austria



universität  
wien



## Distributed algorithms for matrix computations

- Decentralized
  - Nodes operate (mostly) with local information
  - Nodes do not need to be synchronized
  - Automatically adapt to arbitrary topologies (incl. changes during runtime)
- ⇒ Less synchronization, less global information than classical parallel algorithms
- ⇒ Dynamically changing communication schedules

## Which potential do they have in terms of

- Scalability with system size ?
- Resilience ?

## Distributed Data Aggregation Algorithms (DDAAs)

- Reduction operations (summation, averaging, etc.)
- Based on (randomized) **gossiping protocols**:

- 
- 1: **loop** (in node  $k$ )
  - 2:   send data to randomly chosen neighbor
  - 3:   if received data  $\Rightarrow$  update local data
  - 4: **end loop**
- 

- Well established in distributed systems, sensor networks, telecommunications, etc.

Three levels:

- DDAAAs
- Distributed BLAS operations
- Matrix computations

Case study in this talk: orthogonalization / QR decomposition

- 1 Introduction
- 2 Distributed Data Aggregation
  - Existing Methods
  - Resilience
  - The Push-Flow Algorithm
  - Numerical Experiments
- 3 Robust Distributed mGS
- 4 Discussion
- 5 Conclusions and Outlook

## Push-sum

[Kempe et al. 2003]

- Send half of the local value  $x_i$  to a randomly chosen node
- Update a weight  $w_i$  such that  $x_i/w_i$  is the local estimate
- Guaranteed to converge linearly to sum or average

## LiMoSense

[Eyal et al. 2011]

- Push-sum + history

## Flow updating

[Jesus et al. 2009]

- Maintain a *flow variable* for each communication link
  - Local value is added to the flow variable and then the flow variable is sent to randomly chosen node
  - Receiver updates its own flow variable with the negated received flow
- ⇒ Without failures, sum of flows is zero (*flow conservation*, cf. network flow algorithms)
- Recovering from a failure corresponds to (re)establishing flow conservation
  - Local estimate = subtract the sum of flows maintained from initial value, then average with all local estimates of neighbors
  - Convergence (speed) not formally analyzed (appears slow)

- (F1) *Reported temporary* unavailability of nodes/links
- (F2) *Unreported* loss or corruption of a message
- (F3) *Reported permanent* node/link failures
- (F4) *Unreported* corruption of data (e. g., bit flip)
- (F5) *Unreported permanent* node failures



# Resilience Properties of DDAAAs

	Push-sum	LiMoSense	Flow updating
(F1)	✓	✓	✓
(F2)	—	✓	✓
(F3)	—	✓	✓
(F4)	—	—	✓
(F5)	—	—	—

## Note:

Flow-based approach can also recover from purely local failures of a node (F4)

## ⇒ Idea:

- Integrate the flow concept into push-sum !

## Benefits:

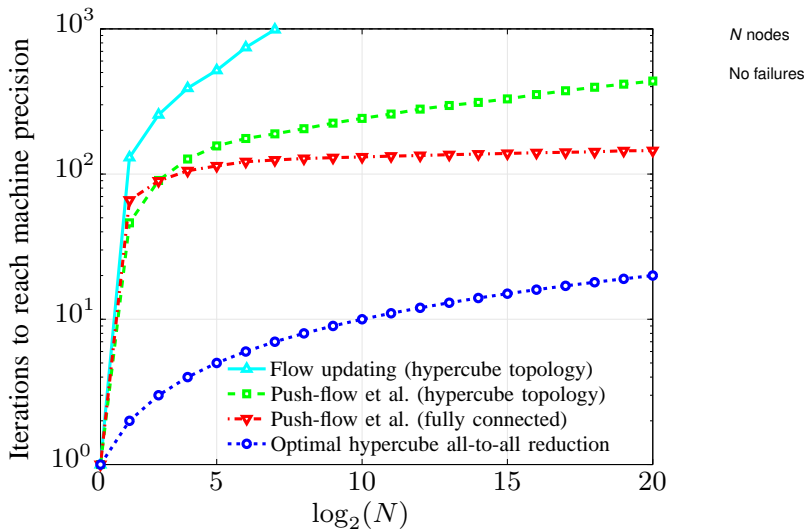
- Convergence properties of push-sum
- Improved resilience due to flow concept

# The Push-Flow Algorithm

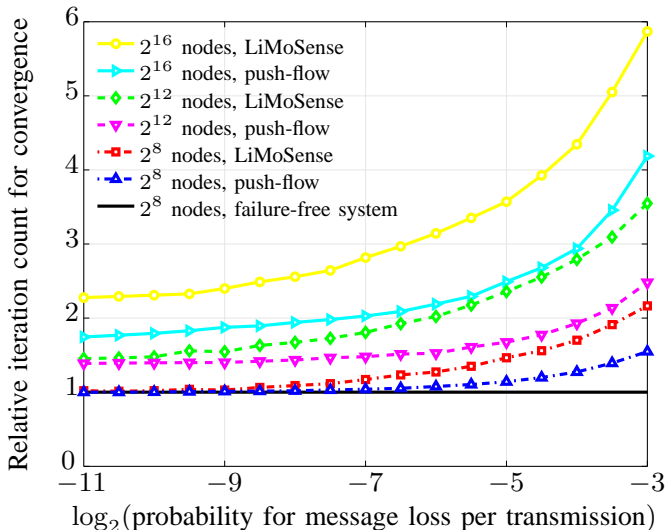
- 
- 1: initialize:  $v_i \leftarrow (x_i, 1)$ ,  $f_{i,j} \leftarrow (0,0)$
  - 2: **for all** received pairs  $f_{j,i}$  **do**
  - 3:      $f_{i,j} \leftarrow -f_{j,i}$
  - 4: **end for**
  - 5: choose a random neighbor  $k \in \mathcal{G}_i$
  - 6: update the flow to node  $k$ :  $f_{i,k} \leftarrow f_{i,k} + (v_i - \sum_{j \in \mathcal{G}_i} f_{i,j}) / 2$
  - 7: send  $f_{i,k}$  to node  $k$
- 

$\mathcal{G}_i$  denotes node  $i$ 's neighborhood

# Scaling Behavior



# Scaling Behavior



- 1 Introduction
- 2 Distributed Data Aggregation
- 3 Robust Distributed mGS**
  - The Algorithm
  - Numerical Experiments
- 4 Discussion
- 5 Conclusions and Outlook

---

**Input:**  $A \in \mathbb{R}^{n \times m}$  (for simplified illustration  $n=N$ )

**Output:**  $Q \in \mathbb{R}^{n \times m}$ ,  $R \in \mathbb{R}^{m \times m}$

```
1: for  $i = 1$  to  $m$  do (in node  $k$ )
2:
3:    $x(k) = A(k, i)^2$ 
4:    $s = \sum_{l=1}^n x(l)$ 
5:    $R(i, i) = \sqrt{s}$ 
6:
7:    $Q(k, i) = A(k, i) / R(i, i)$ 
8:
9:
10:  for  $j = i + 1$  to  $m$  do
11:
12:     $x(k) = Q(k, i)A(k, j)$ 
13:     $R(i, j) = \sum_{l=1}^n x(l)$ 
14:
15:     $A(k, j) = A(k, j) - Q(k, i)R(i, j)$ 
16:
17:
18:  end for
19: end for
```

---

---

**Input:**  $A \in \mathbb{R}^{n \times m}$  (for simplified illustration  $n=N$ )

**Output:**  $Q \in \mathbb{R}^{n \times m}$ ,  $R \in \mathbb{R}^{m \times m}$

```
1: for  $i = 1$  to  $m$  do (in node  $k$ )
2:
3:    $x(k) = A(k, i)^2$ 
4:    $s_k = \text{DDAA}(x)$ 
5:    $R_k(i, i) = \sqrt{s_k}$ 
6:
7:    $Q(k, i) = A(k, i) / R_k(i, i)$ 
8:
9:
10:  for  $j = i + 1$  to  $m$  do
11:
12:     $x(k) = Q(k, i)A(k, j)$ 
13:     $R_k(i, j) = \text{DDAA}(x)$ 
14:
15:     $A(k, j) = A(k, j) - Q(k, i)R_k(i, j)$ 
16:
17:
18:  end for
19: end for
```

---



---

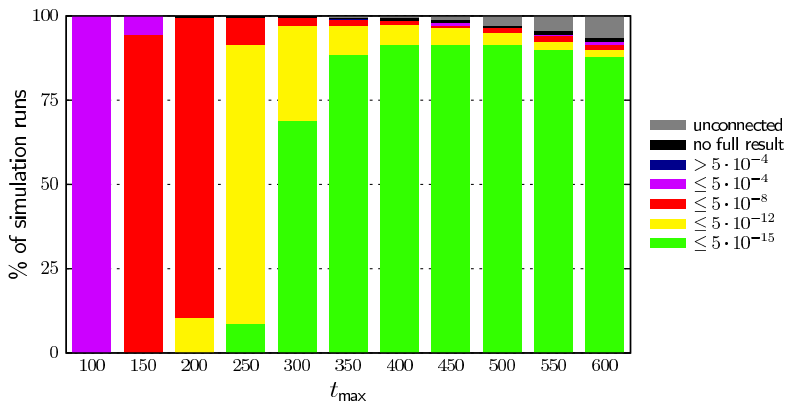
**Input:**  $A \in \mathbb{R}^{n \times m}$  (for simplified illustration  $n=N$ )

**Output:**  $Q \in \mathbb{R}^{n \times m}$ ,  $R \in \mathbb{R}^{m \times m}$

```
1: for  $i = 1$  to  $m$  do (in node  $k$ )
2:   ... check for node failures, update  $P_k$  and  $B_k$  ...
3:    $x(k) = \sum_{p \in P_k} A(p, i)^2$ 
4:    $s_k = \text{DDAA}(x)$ 
5:    $R_k(i, i) = \sqrt{s_k}$ 
6:   for each  $p \in P_k$ 
7:      $Q(p, i) = A(p, i) / R_k(i, i)$ 
8:   for each  $b \in B_k$ 
9:      $Q(b, i) = A(b, i) / R_k(i, i)$ 
10:  for  $j = i + 1$  to  $m$  do
11:    ... check for node failures, update  $P_k$  and  $B_k$  ...
12:     $x(k) = \sum_{p \in P_k} Q(p, i) A(p, j)$ 
13:     $R_k(i, j) = \text{DDAA}(x)$ 
14:    for each  $p \in P_k$ 
15:       $A(p, j) = A(p, j) - Q(p, i) R_k(i, j)$ 
16:    for each  $b \in B_k$ 
17:       $A(b, j) = A(b, j) - Q(b, i) R_k(i, j)$ 
18:  end for
19: end for
```

---

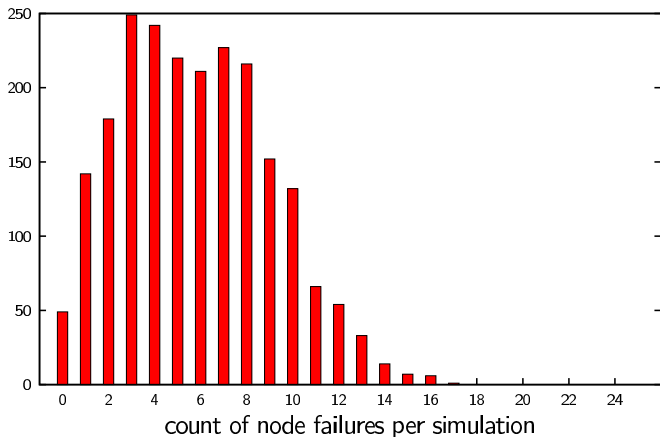
# Factorization Error of rdmGS



5D hypercube,  $\lambda = 12$  [s], 200 simulation runs

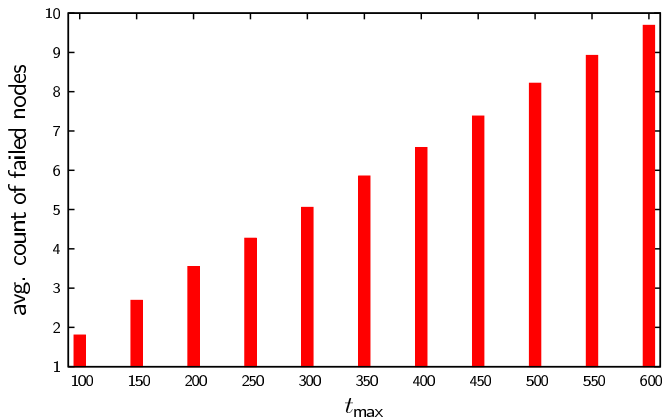
On average, 5.82 nodes failed per simulation run (min = 0, max = 17)

# Node Failures in Simulation



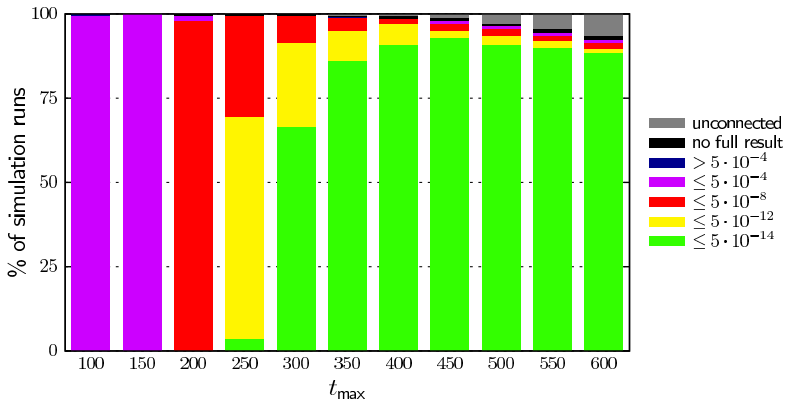
**Figure:** Number of node failures per simulation for  $\lambda = 12[s]$  and  $h = 5$  over all  $t_{\max}$ . Average of 5.82 node failures over all simulations

# Node Failures in Simulation



**Figure:** Average number of node failures per simulation for  $\lambda = 12$ [s] and  $h = 5$  and varying  $t_{\max}$ . Average of 5.82 node failures over all simulations

# Orthogonality of rdmGS



5D hypercube,  $\lambda = 12$  [s], 200 simulation runs

On average, 5.82 nodes failed per simulation run (min = 0, max = 17)

- 1 Introduction
- 2 Distributed Data Aggregation
- 3 Robust Distributed mGS
- 4 Discussion**  
Alternative Approaches to Resilience
- 5 Conclusions and Outlook

## Influence of the context

Performance (comparison) of various approaches strongly depends on various context parameters !

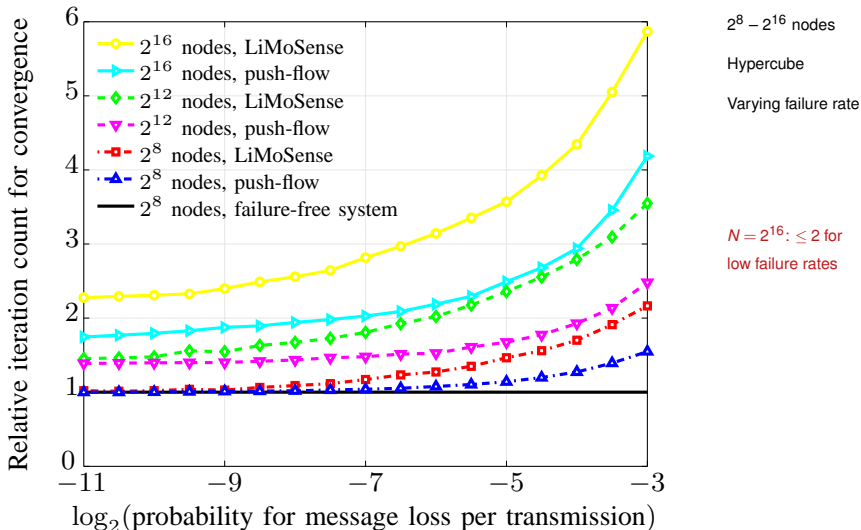
- Topology
- Routing information (beyond local neighborhood)
- Properties of the failure distribution(s):  
MTBF, MTTR, etc.
- Properties of the application:  
ratio of application runtime to checkpointing interval, etc.
- ...

# Checkpointing & Restarting

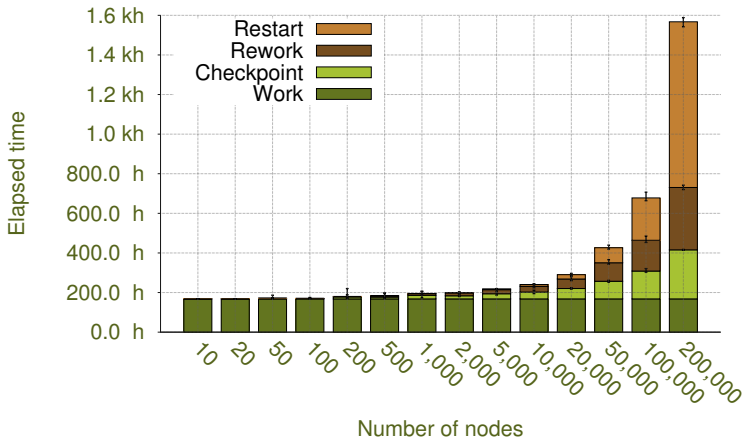
- Overhead in terms of time and storage
- Often assumes restart of the application on the same number of nodes
- Coordinated vs. uncoordinated
- Stable common storage vs. distributed storage



# Overhead of DDAs



# Overhead of Coordinated Checkpointing



Wall clock times for an application based on coordinated checkpointing for increasing number of nodes

Source: [Varela, Ferreira and Riesen; 2010]

## Redundant computing

e.g., [Engelmann et al. 2009, Ferreira et al. 2008]

- Spare nodes are usually required (avoid imbalance)  
⇒ large hardware overhead
- rdmGS integrates redundant computing and (reactive) migration concepts *without extra hardware*

## Algorithm-based fault tolerance (ABFT)

e.g., [Huang & Abraham 1984, Chen & Dongarra 2008, Chen 2011]

- Extend input by checksums, detect and recover from errors
  - Usually at a higher level than elementary data aggregation
  - Deterministic correction vs. randomized “healing”
  - Communication overhead vs. slow-down of convergence
- Could complement each other ?

## Current status

- DDAAAs: “self-healing” methods
  - Lower overhead
  - Scale asymptotically like parallel reduction
  - Failures slow down convergence
- rdMGS: redundancy on top of DDAAAs
- Concept has potential
- Details depend strongly on context parameters
- Performance penalty in practice to be investigated

## Central questions

- Runtime performance comparison ?
- Depends on
  - Topology
  - Failure rates (ratio to total runtime)
  - (Granularity and structure of) Application
  - Checkpointing interval (ratio to total runtime of application)
  - ...

## Work in progress

- Quantitative performance evaluation and comparison
- Convergence acceleration
  - The more global information (topology, routing, etc.) you can utilize, the faster you can make it !

*Thank you for your attention !*

`http://rlcta.univie.ac.at`