

# Verification of Resilient Communication Models for the Simulation of a Highly Adaptive Energy-Efficient Computer

Mario Bielert<sup>§</sup>, Kim Feldhoff<sup>§</sup>, Florina M. Ciorba<sup>‡</sup>, Stefan Pfennig<sup>§</sup>, Elke Franz<sup>§</sup>, Thomas Ilsche<sup>§</sup>, and Wolfgang E. Nagel<sup>§</sup>

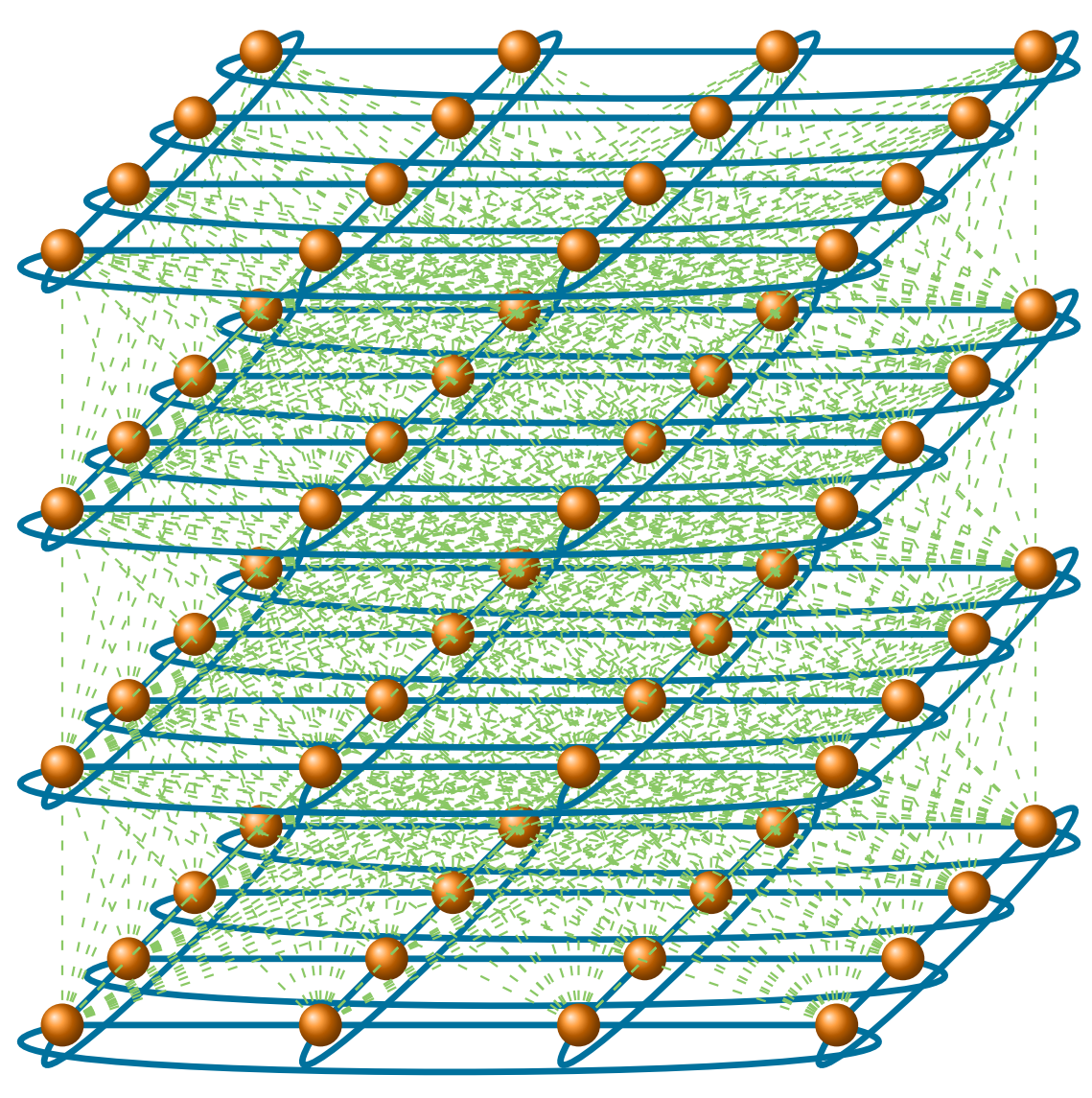
## Motivation and Goal

- **HAEC Box:** high performance–low energy computing box [1].
  - Optical (on-board) chip-to-chip communication.
  - Wireless (intra-board) chip-to-chip communication.
  - Runtime software adaptivity.
- ⇒ Hardware adapts to the needs of the computational problem.
- **HAEC-SIM:** integrated simulation environment for studying the performance and energy costs of the HAEC Box [2].
  - Simulator design based on individual abstraction models of
    - \* Hardware (e.g., CPUs, links),
    - \* Architecture (e.g., computing nodes, network),
    - \* Software (e.g., runtime system, code generation).

• **Goal:** The goal of this work is the verification of the resilient communication models [3] implemented in HAEC-SIM [2].

### HAEC Box specifications

- 4 boards, each with 4 × 4 computing nodes.
- Each computing node comprises 3D stacked processor chips which contain a large number of “thin” cores.



● computing node — optical link — wireless link

Table 1: Characteristics of the optical and wireless links of the HAEC Box.

	Optical links	Wireless links
Topology	2D torus	fully connected crossbar
Bandwidth	250 Gbit/s	100 Gbit/s
Latency	10 ns	100 ns
Bit error rate	10 <sup>-12</sup>	10 <sup>-8</sup>

## Applications

- The benchmarks LU.D.4096 and BT.D.4096 of the NAS Parallel Benchmark Suite 3.3 are chosen as example applications.
- Both benchmarks are communication intensive: LU.D.4096 spends approximately 68 % in MPI functions, while BT.D.4096 spends approximately 48 % (Table 2).

- The communication matrices of both benchmarks are not dispersed: Processes communicate primarily with their neighboring ranks.
- 99 % of messages exchanged in LU.D.4096 are of size 240 B and 280 B, respectively.
- In BT.D.4096, most exchanged messages are of sizes 11.484 KiB and 1.914 KiB, while other larger messages of size 200 MiB represent at most 1 % of the exchanged point-to-point messages.
- The 4096 MPI processes of the benchmarks have been executed on 256 Taurus [4] nodes, as 16 processes per node.

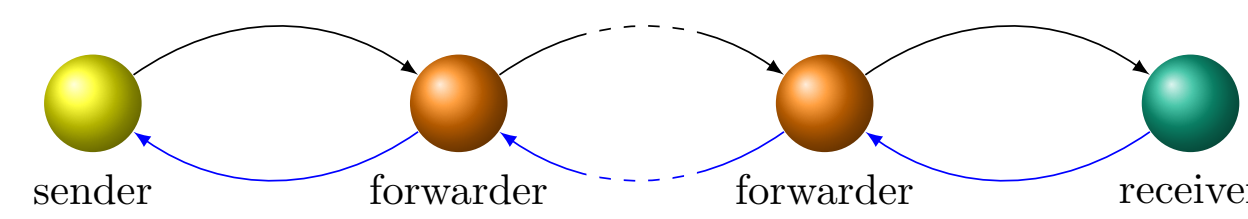
Table 2: Total running times and percentage of exclusive accumulated times for Application and MPI functions on Taurus.

Benchmark	Running time	Exclusive Application %	Exclusive MPI %
LU.D.4096	19.625 s	32.40 %	67.60 %
BT.D.4096	24.710 s	52.07 %	47.93 %

## Simulated HPC System Specifications

- 3D torus link topology, with 16 × 16 × 16 computing nodes.
  - This topology is a intermediate step between current typical HPC systems and the HAEC Box topology.
- Communication links
  - Infiniband: 700 ns latency and 54,54 Gbit/s bandwidth.
  - Assumed to be congestion free, as the implementation of resource contention protocols is ongoing.
- Packet loss probability per link: 0.01.
- Mapping: During simulation, the MPI processes of the benchmarks are mapped to the simulated compute nodes in an xyz order.

## Unicast (or Point-to-point) Communication



→ sending a digitally signed data packet  
← sending a digitally signed acknowledgment

- Sender splits messages into multiple data packets of equal length (1500 B).
- Data packets and acknowledgments are digitally signed by sender and receiver to prevent unrecognized modifications.

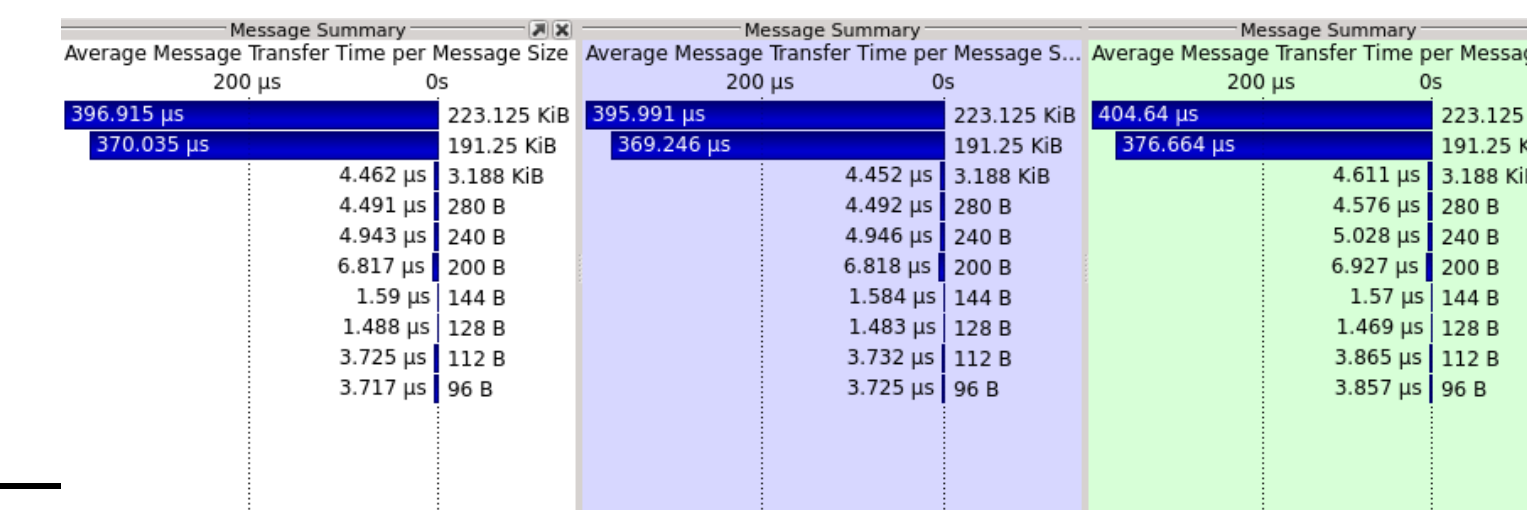
## Resilient Communication Models

**DOR, (resilient dimension order routing):** common store-and-forward approach of sending packets; packets are organized in windows; acknowledgments for single data packets.

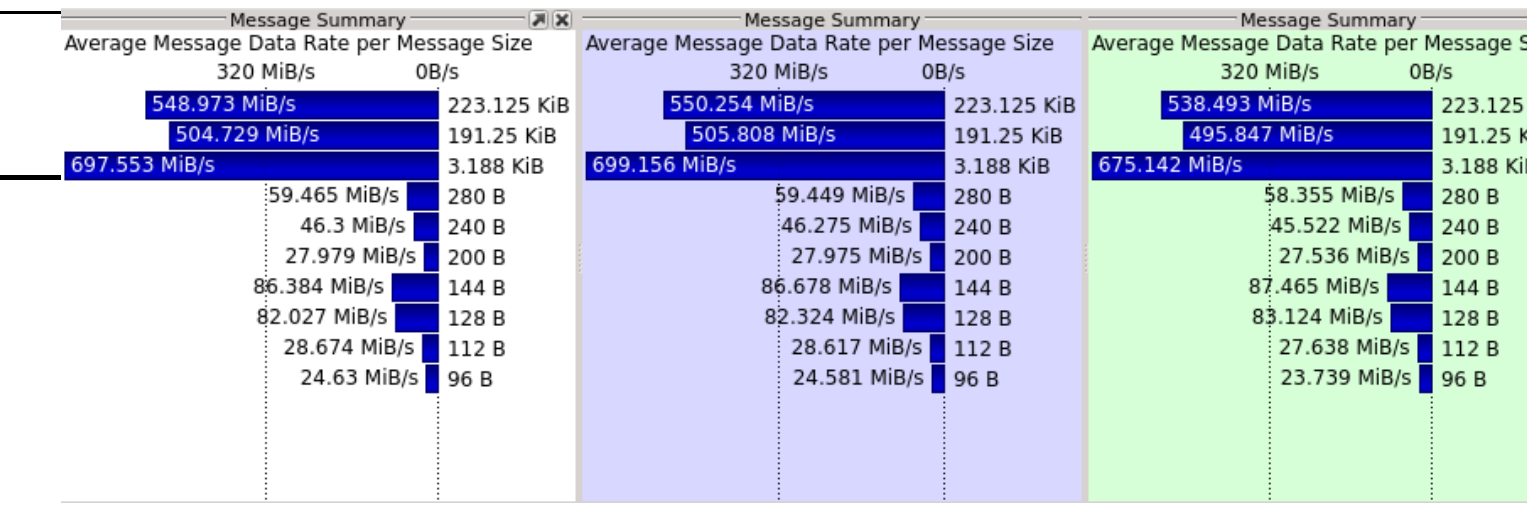
**PNC, (resilient practical network coding):** sender computes and sends random linear combinations out of the native packets of one generation; acknowledgments confirm current rank of the matrix of received packets.

**NCD, (resilient dynamical network coding):** similar to PNC, but the sender node computes the number of packets to be sent based on an estimation of the delivery probability [5].

## Message Statistics

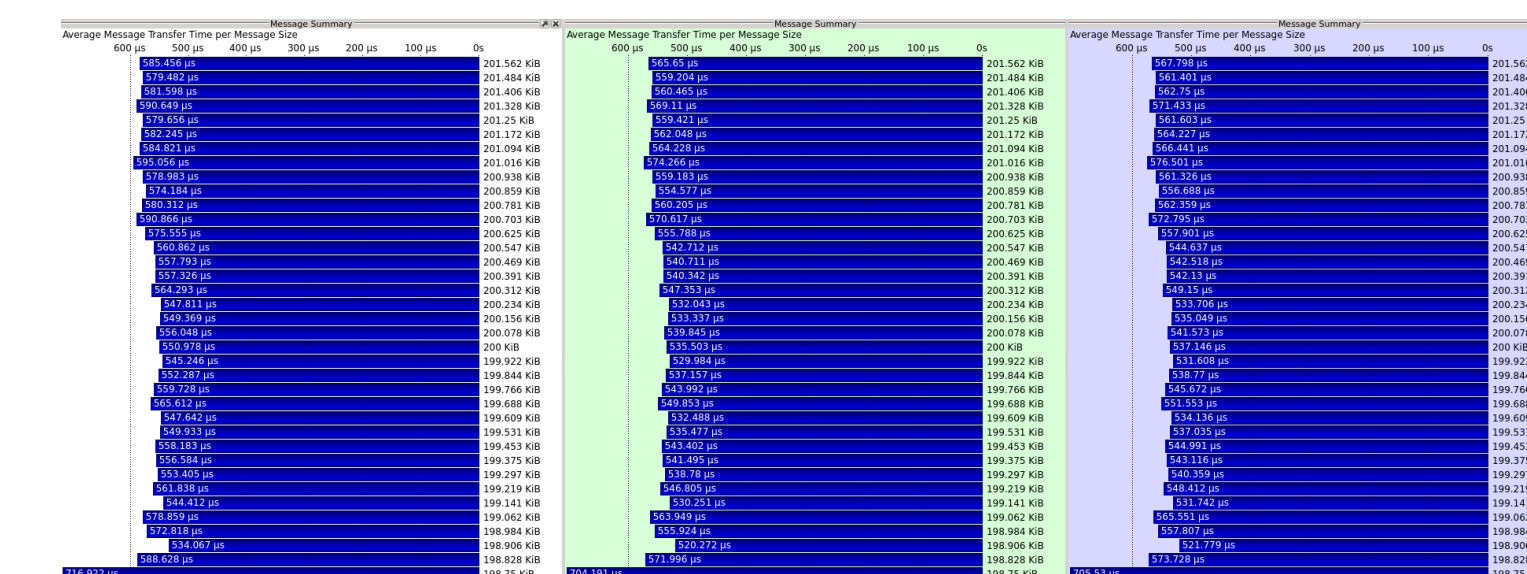


(a) LU.D.4096 – Message transfer times per message size

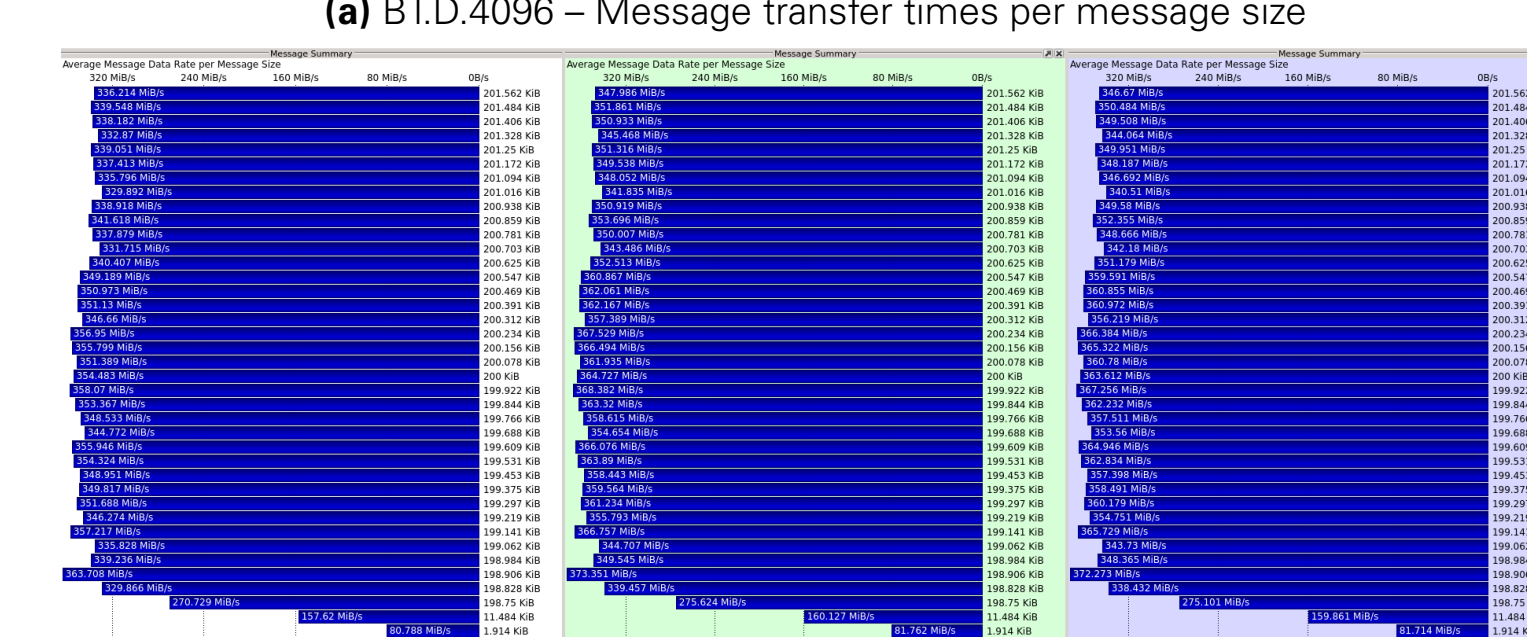


(b) LU.D.4096 – Message data rates per message size

Figure 1: (a) Message transfer times per message size, and (b) Message data rates per message size, for the simulated execution of LU.D.4096 on the target platform using PNC, (left), NCD, (middle), and DOR, (right) as communication models.



(a) BT.D.4096 – Message transfer times per message size



(b) BT.D.4096 – Message data rates per message size

Figure 2: (a) Message transfer times per message size, and (b) Message data rates per message size, for the simulated execution of BT.D.4096 on the target platform using PNC, (left), NCD, (middle), and DOR, (right) as communication models.

- The message data rates (Figures 1 and 2) are comparable yet clearly distinguishable for each communication model:

- PNC<sub>r</sub> and DOR<sub>r</sub> results in the lowest data rates for BT and LU, respectively.
- NCD<sub>r</sub> achieves the highest data rates.

## Communication-to-computation Ratio

The differences in running times and MPI-communication to application-computation ratios for LU.D.4096 (Table 3) are small, yet PNC<sub>r</sub> shows the best performance, very closely followed by NCD<sub>r</sub>.

Table 3: Total running times and percentage of exclusive accumulated times for LU.D.4096 for Application and MPI functions.

Model	Running time	Exclusive Application %	Exclusive MPI %
PNC <sub>r</sub>	23.723 s	26.44 %	73.56 %
NCD <sub>r</sub>	23.727 s	26.43 %	73.57 %
DOR <sub>r</sub>	23.827 s	26.32 %	73.68 %

The differences in running times and MPI communication-to-Application computation ratios for BT.D.4096 (Table 4) are also small, yet NCD<sub>r</sub> shows the best performance, very closely followed by PNC<sub>r</sub>.

Table 4: Total running times and percentage of exclusive accumulated times for BT.D.4096 for Application and MPI functions.

Model	Running time	Exclusive Application %	Exclusive MPI %
PNC <sub>r</sub>	24.847 s	47.62 %	52.38 %
NCD <sub>r</sub>	24.840 s	47.63 %	52.37 %
DOR <sub>r</sub>	24.913 s	47.47 %	52.53 %

**Remark:** The values of Tables 3 and 4 cannot be meaningfully compared to those in Table 2. This is due to a changed topology, oversubscription of nodes and different error-probability rates between the original platform and the target platform.

## Number of Hops per Message

Distribution of the number of hops per message

- For the target HPC platform topology.
- Depends on processes-to-nodes mapping and communication patterns.
- Influences communication-to-computation ratio.

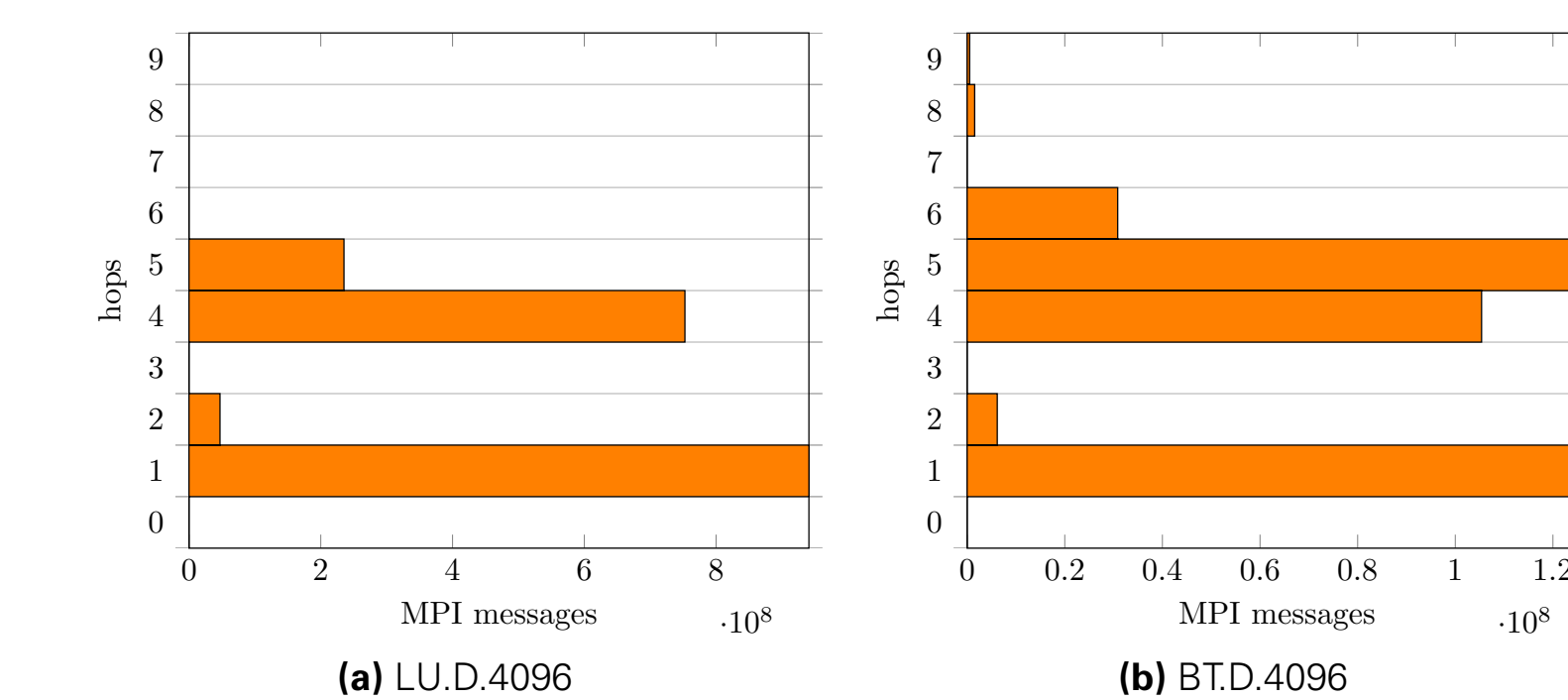


Figure 3: Number of hops travelled by the point-to-point messages for the two applications considered on the target platform topology (3D torus).

## Determination of the Message Transfer Times

The message transfer times are calculated by means of a polynomial approximation in HAEC-SIM. Figure 4 shows no significant

error between simulation results obtained with Sage and the polynomial approximation of HAEC-SIM (e.g., 0.248 % error for 6 hops).

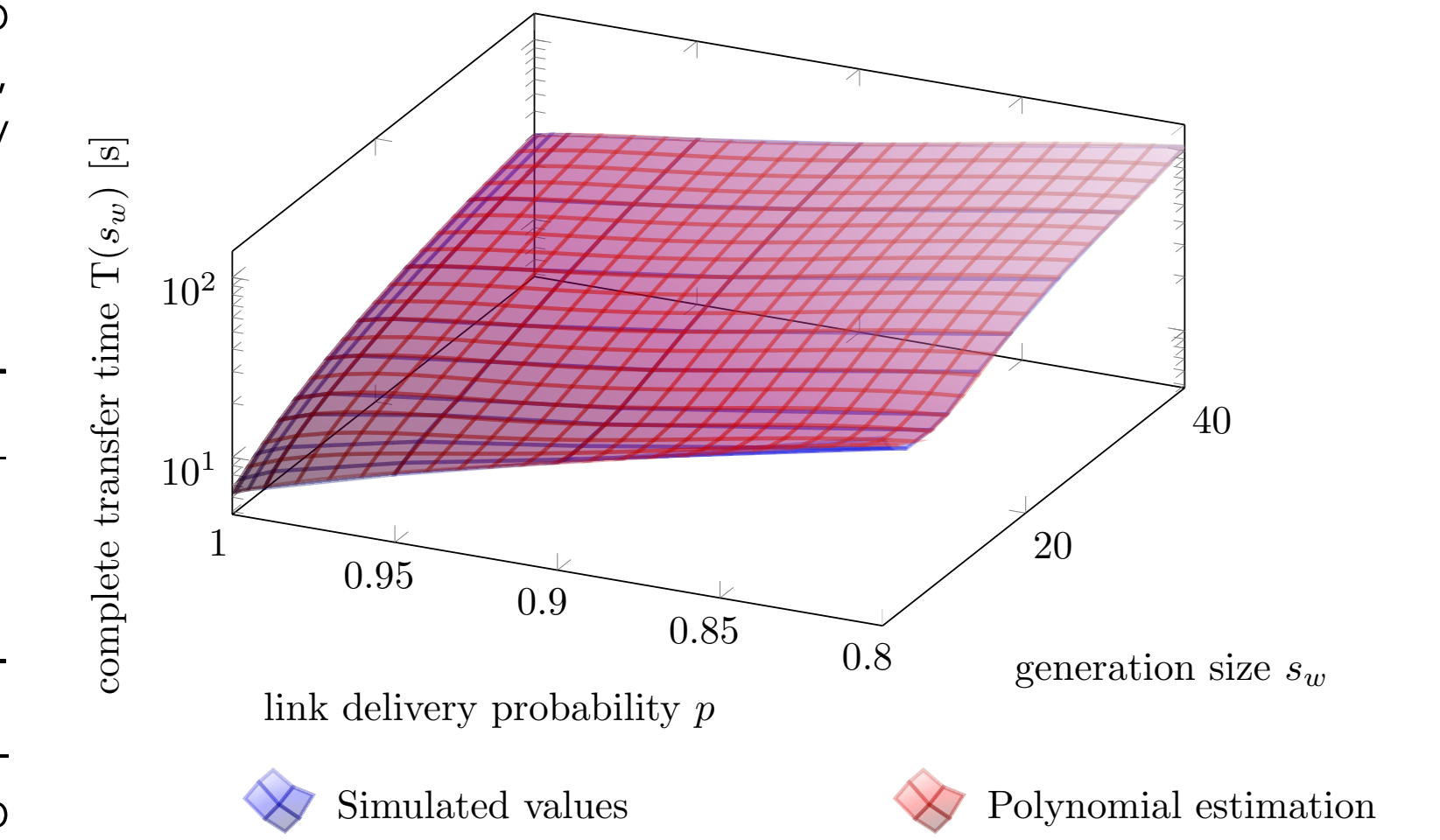


Figure 4: Comparison of simulated values to polynomial approximation for 6 hops.

## Verification

- The workflow for the verification of the implementation of the resilient network models in HAEC-SIM (Figure 5) comprises
  - Polynomial approximations of the simulation values obtained with Sage [3] and
  - Comparison of results of the independent polynomial-based implementations in HAEC-SIM and in the verification tool, respectively.

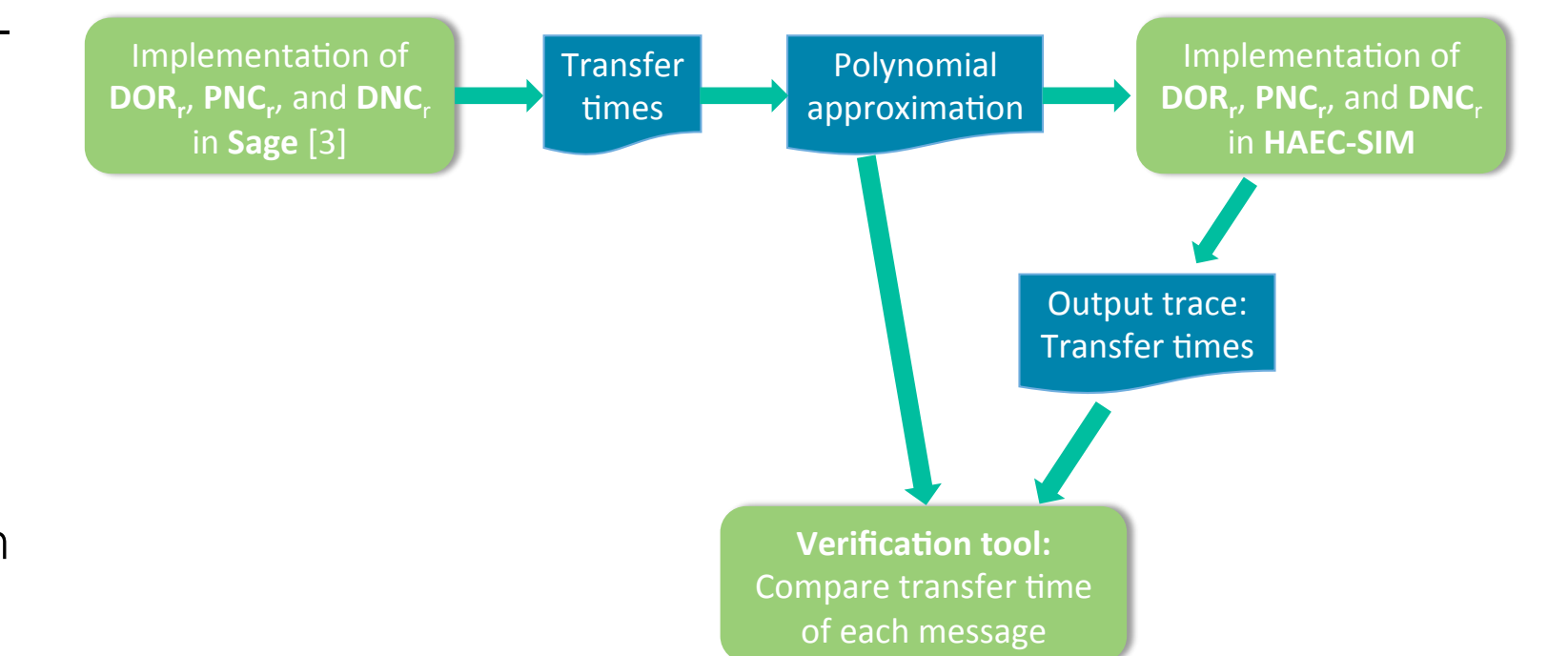


Figure 5: Workflow for the verification of the implementation in HAEC-SIM.

- The verification showed, that each message simulated with HAEC-SIM conforms to the independent implementation in the verification tool.

## Conclusions and Future Work

- This work represents the first step towards accurate HAEC-SIM-based simulations for studying the behavior of communication intensive applications running on the HAEC Box using three resilient communication models.
- Verification of the simulated results against those from the independent implementation shows consistency.
- Of the investigated models, PNC<sub>r</sub> and NCD<sub>r</sub> perform similar, and both outperform DOR<sub>r</sub>.
- Future work directions include developing communication models that consider link congestion and collective communication.

## References

- [1] G. Fettweis, W. Nagel, and W. Lehner, “Pathways to servers of the future,” in *Proc. of the Design, Automation Test in Europe Conference Exhibition*, Mar 2012.
- [2] M. Bielert, F. M. Ciorba, K. Feldhoff, T. Ilsche, and W. E. Nagel, “HAEC-SIM: A Simulation Framework for Highly Adaptive Energy-Efficient Computing Platforms,” in *Proc. of the Conference on Simulation Tools and Techniques*, accepted, 2015.
- [3] S. Pfennig, E. Franz, F. M. Ciorba, T. Ilsche, and W. E. Nagel, “Modeling communication delays for network coding and routing for error-prone transmission,” in *Proc. of the 3rd Intl. Conf. on Future Gen. Comm. Techn.* IEEE, Aug 2014, pp. 19–24.
- [4] Technische Universität Dresden, ZIH, “Cluster Taurus2,” <https://doc.zih.tu-dresden.de/hpc-wiki/bin/view/Compendium/HardwareTaurus>.
- [5] S. Pfennig and E. Franz, “Adjustable Redundancy for Secure Network Coding in a Unicast Scenario,” in *Proc. of International Symposium on Network Coding*, 2014.