

Flow-level Tail Latency Estimation and Verification based on Extreme Value Theory

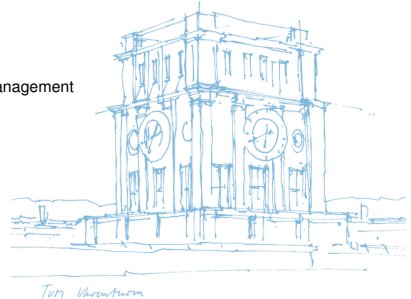
Max Helm, Florian Wiedner, and Georg Carle

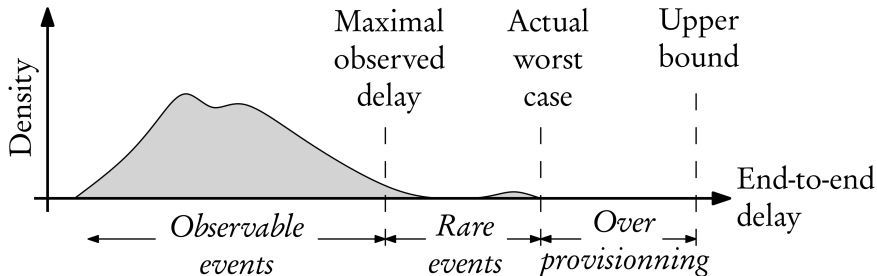
November 1, 2022

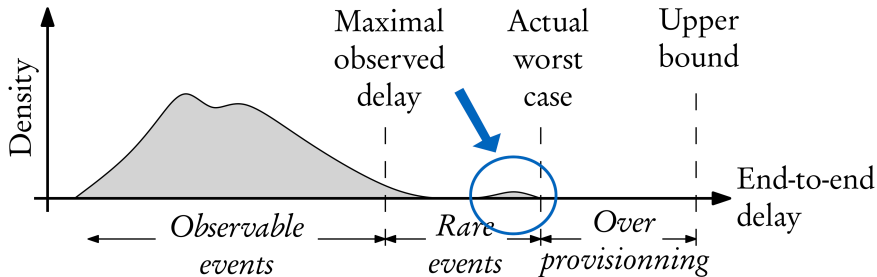
18th International Conference on Network and Service Management

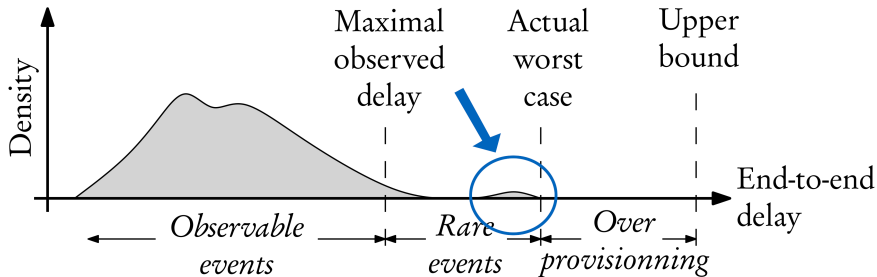
Thessaloniki, Greece

Chair of Network Architectures and Services
Department of Informatics
Technical University of Munich

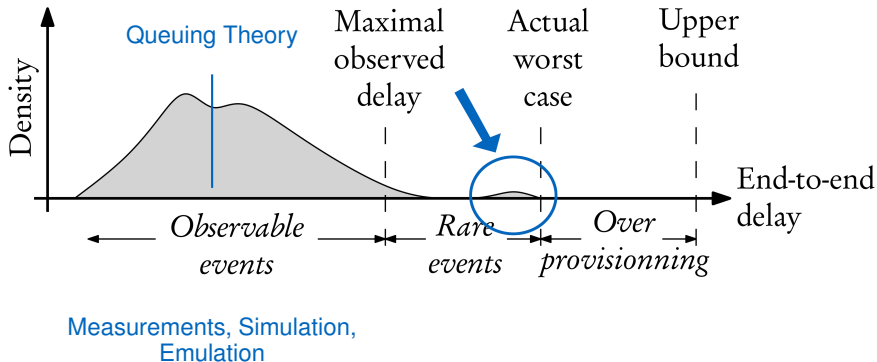


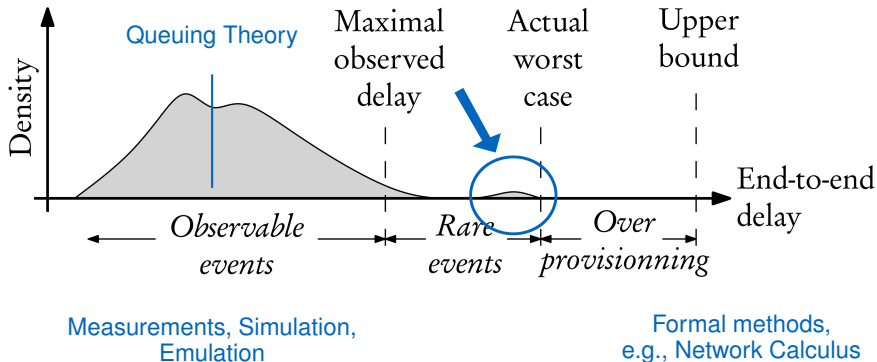


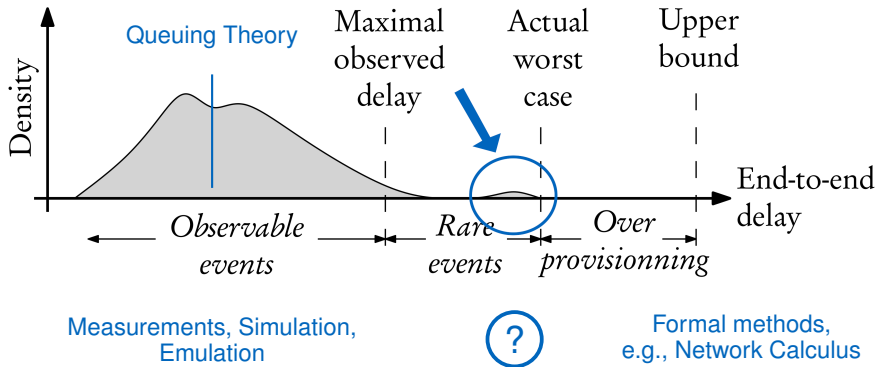


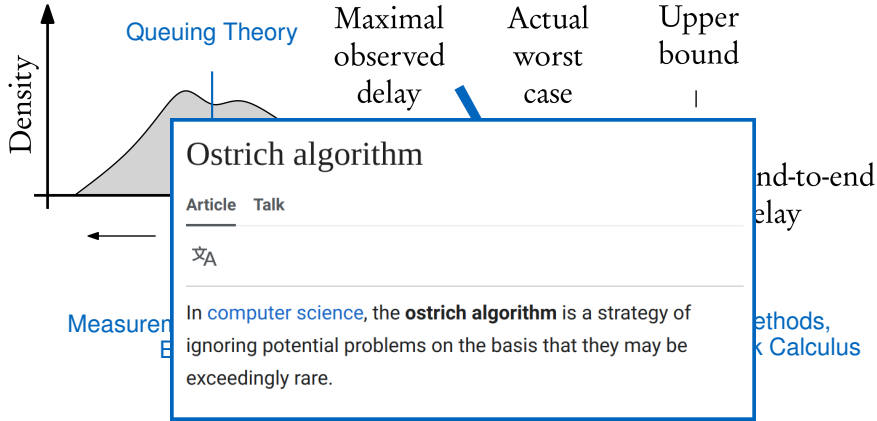


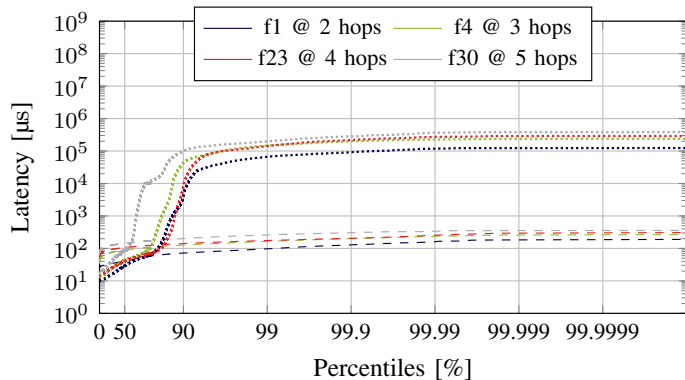
Measurements, Simulation,
Emulation











- End-to-end latency measurements of multihop flows¹
- Different flow lengths and measurement types (emulation, virtualized hardware measurements)
- Extreme latency spikes at different percentiles

¹ Wiedner, Florian, et al. "HVNet: Hardware-Assisted Virtual Networking on a Single Physical Host." INFOCOM WKSHPS CNERT 2022.

- Expected latency behavior at high quantiles, e.g., 99.999th percentile
- **Magnitude of rare events** or frequency of events of given magnitude
- Based on as few data points as possible, i.e., requiring only short measurement periods
- Fast calculation, avoiding computationally expensive simulations, emulations, or time consuming measurements

- Expected latency behavior at high quantiles, e.g., 99.999th percentile
- **Magnitude of rare events** or frequency of events of given magnitude
- Based on as few data points as possible, i.e., requiring only short measurement periods
- Fast calculation, avoiding computationally expensive simulations, emulations, or time consuming measurements

Questions we want to answer:

Given a few latency measurements per flow in a network:

- Service Level Agreement clause: Given latency is exceeded only once during given time period
- Is the latency behavior **bounded or unbounded**?
- What value is the latency **converging** to?

Extreme Value Theory (EVT):

“Extreme value theory is unique as a statistical discipline in that it develops techniques and models for describing the unusual rather than the usual.”

— Coles, Stuart, et al. An Introduction to Statistical Modeling of Extreme Values. Vol. 208. London: Springer, 2001.

- Commonly used to **predict rare events** such as storms or floods
- Models the tail of distributions
- Model can be used to predict occurrence of rare events belonging to the tail of the distribution

Steps to obtain an EVT model:

1. Select a **threshold**, indicating which values belong to the tail
2. Fit all values above the threshold to a **Generalized Pareto Distribution** (GPD)
3. GPD is defined by three parameters: Threshold (μ), Location (σ), and Tail (ξ)

Steps to obtain an EVT model:

1. Select a **threshold**, indicating which values belong to the tail
2. Fit all values above the threshold to a **Generalized Pareto Distribution** (GPD)
3. GPD is defined by three parameters: Threshold (μ), Location (σ), and Tail (ξ)

Steps to evaluate an EVT model:

- Predict occurrence of events using the GPD, check if they match observations
- Can be achieved using the **Return Level**
- Return Level is the value that is expected to be exceeded on average exactly once during a given Return Period

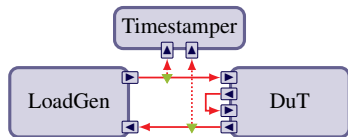
or

- **Compare quantiles** of EVT model to empirical quantiles of evaluation data

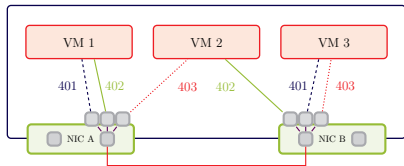
Methodology

Latency Measurements²

Hardware setup:

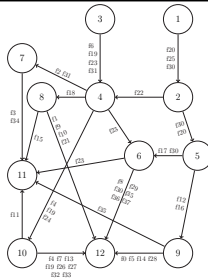


Virtualized nodes:



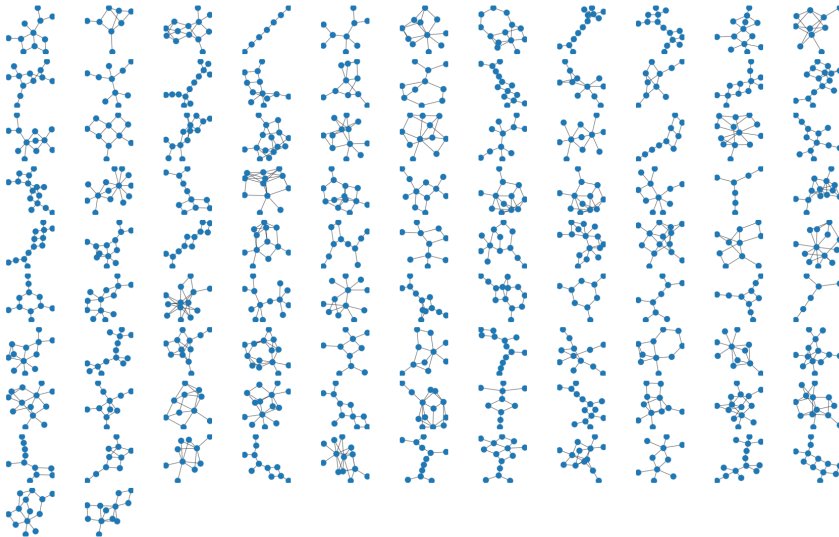
100 random network topologies and flow configurations:

Parameter	Minimum	Maximum	Mean	Σ
Number of Network Nodes	6	15	12	1,190
Number of Flows	19	59	35	3,559
Flow Lengths	2	9	3	—
Flow Rates [Mbit/s]	1.0	831	44	—
Link Rates [Mbit/s]	434	2000	705	—
Link Utilization Rates [%]	0	87	24	—

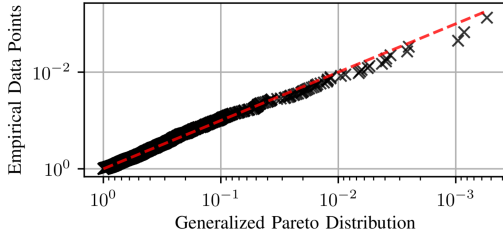


- Per flow latencies
- Total of **14 billion** latency values as input to EVT models

² Wiedner, Florian, et al. "HVNet: Hardware-Assisted Virtual Networking on a Single Physical Host." INFOCOM WKSHPs CNERT 2022.

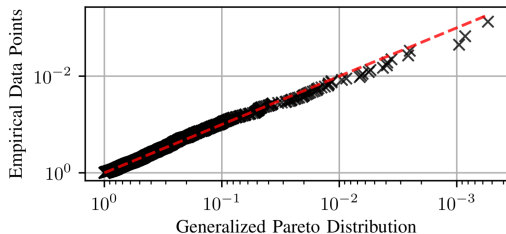


Goodness-of-fit for a Maximum Likelihood Estimator to a GPD:

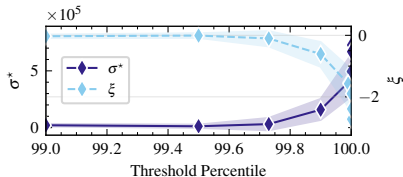


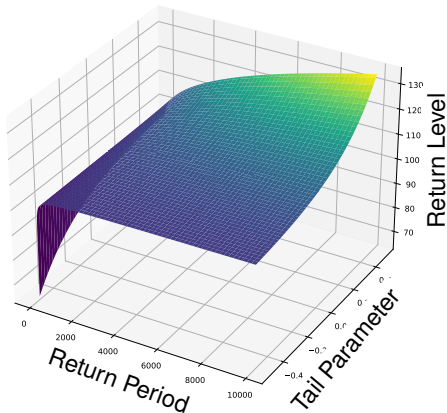
- Generate an EVT model for latencies of each flow
- Maximum Likelihood Estimator (MLE) to fit empirical data points over threshold to GPD
- Threshold selection such that resulting EVT model is stable:

Goodness-of-fit for a Maximum Likelihood Estimator to a GPD:



- Generate an EVT model for latencies of each flow
- Maximum Likelihood Estimator (MLE) to fit empirical data points over threshold to GPD
- Threshold selection such that resulting EVT model is stable:





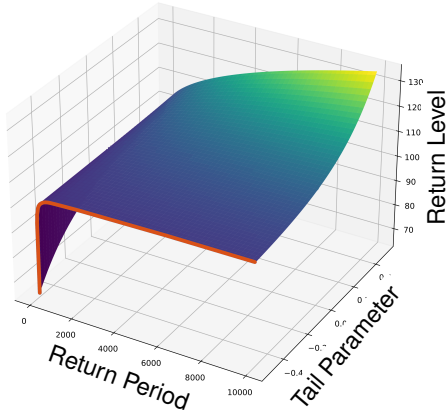
Return Level:

Return level is the value that is on average **exceeded exactly once** during a given return period

$$x_m = \mu + \frac{\sigma}{\xi} \cdot \left[\left(m \cdot \frac{D_{d>\mu}}{D} \right)^\xi - 1 \right]$$

Observations:

- Return level for different values of the tail parameter ξ and the length of the return period m



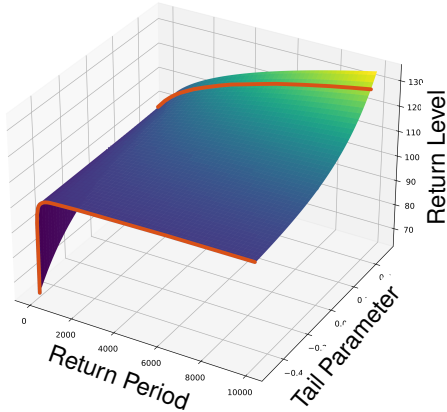
Return Level:

Return level is the value that is on average **exceeded exactly once** during a given return period

$$x_m = \mu + \frac{\sigma}{\xi} \cdot \left[\left(m \cdot \frac{D_{d>\mu}}{D} \right)^\xi - 1 \right]$$

Observations:

- Return level for different values of the tail parameter ξ and the length of the return period m
- $\xi < 0$: Return level **converges** to a fixed value



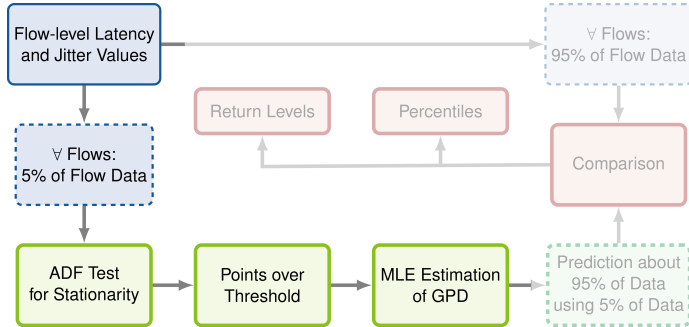
Return Level:

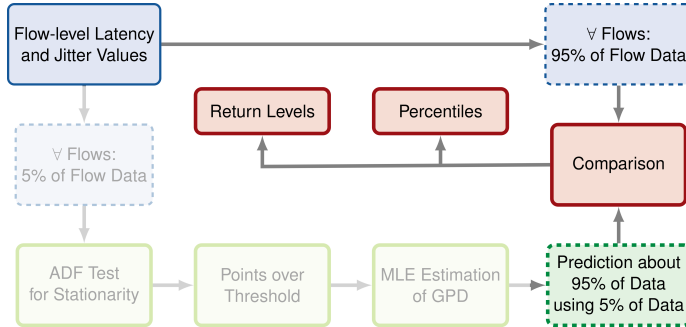
Return level is the value that is on average **exceeded exactly once** during a given return period

$$x_m = \mu + \frac{\sigma}{\xi} \cdot \left[\left(m \cdot \frac{D_{d>\mu}}{D} \right)^\xi - 1 \right]$$

Observations:

- Return level for different values of the tail parameter ξ and the length of the return period m
- $\xi < 0$: Return level **converges** to a fixed value
- $\xi > 0$: Return level **diverges**





Accuracy of return level predictions:

- Return level for 95% of data (unseen), i.e., predictions for a **twentyfold time horizon**
- Return level calculated with confidence intervals of confidence level 95%
- Reducing the time horizon to twofold increases accuracy to 85%

One exceedance	Exceedances $\neq 1$
75%	25%

Evaluation

Return Level

Accuracy of return level predictions:

- Return level for 95% of data (unseen), i.e., predictions for a **twentyfold time horizon**
- Return level calculated with confidence intervals of confidence level 95%
- Reducing the time horizon to twofold increases accuracy to 85%

One exceedance	Exceedances $\neq 1$
75%	25%

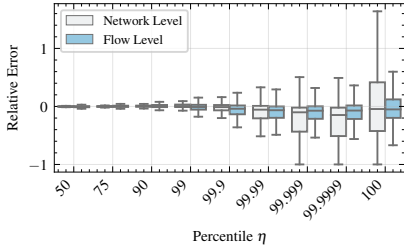
Bounds on return levels:

- Observe bounded as well as un-bounded return levels
- Majority of flows have **bounded return levels**

Bounded Return Level	Unbounded Return Level
3,507 (57.51%)	2,591 (42.49%)

Comparison of percentiles between GPD of EVT model and evaluation data (95% of data points):

Percentile	50	75	90	99	99.9	99.99	99.999	100
MdAPE [%]	0.7	1.0	1.8	4.2	6.8	9.6	11.4	16.8

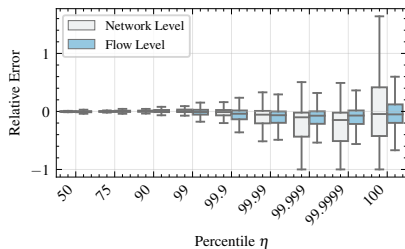


Evaluation

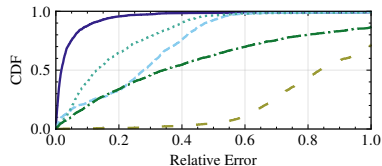
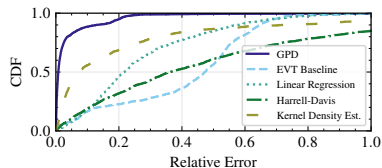
Percentiles

Comparison of percentiles between GPD of EVT model and evaluation data (95% of data points):

Percentile	50	75	90	99	99.9	99.99	99.999	100
MdAPE [%]	0.7	1.0	1.8	4.2	6.8	9.6	11.4	16.8



Comparison to other methods for selected tail percentiles (50th and 90th):



Contributions:

- Flow-level latency EVT models for low-latency virtualized wired networks
- Verification of the approach by testing predictive power of EVT models against twentyfold time periods of unseen latency data
- Comparison of EVT approach against other methods

More details in our paper:

- Related work
- Predictions for latency jitter
- Threshold selection
- Return level accuracy
- Quantile comparisons

Contact:

- helm@net.in.tum.de
- Send me a message on Whova

Link to paper:

