

Methoden und Messverfahren für das automatische Skalieren elastischer Cloud-Umgebungen

Präsentation zur Doktorarbeit

Dr. Nikolas Herbst

Betreuer

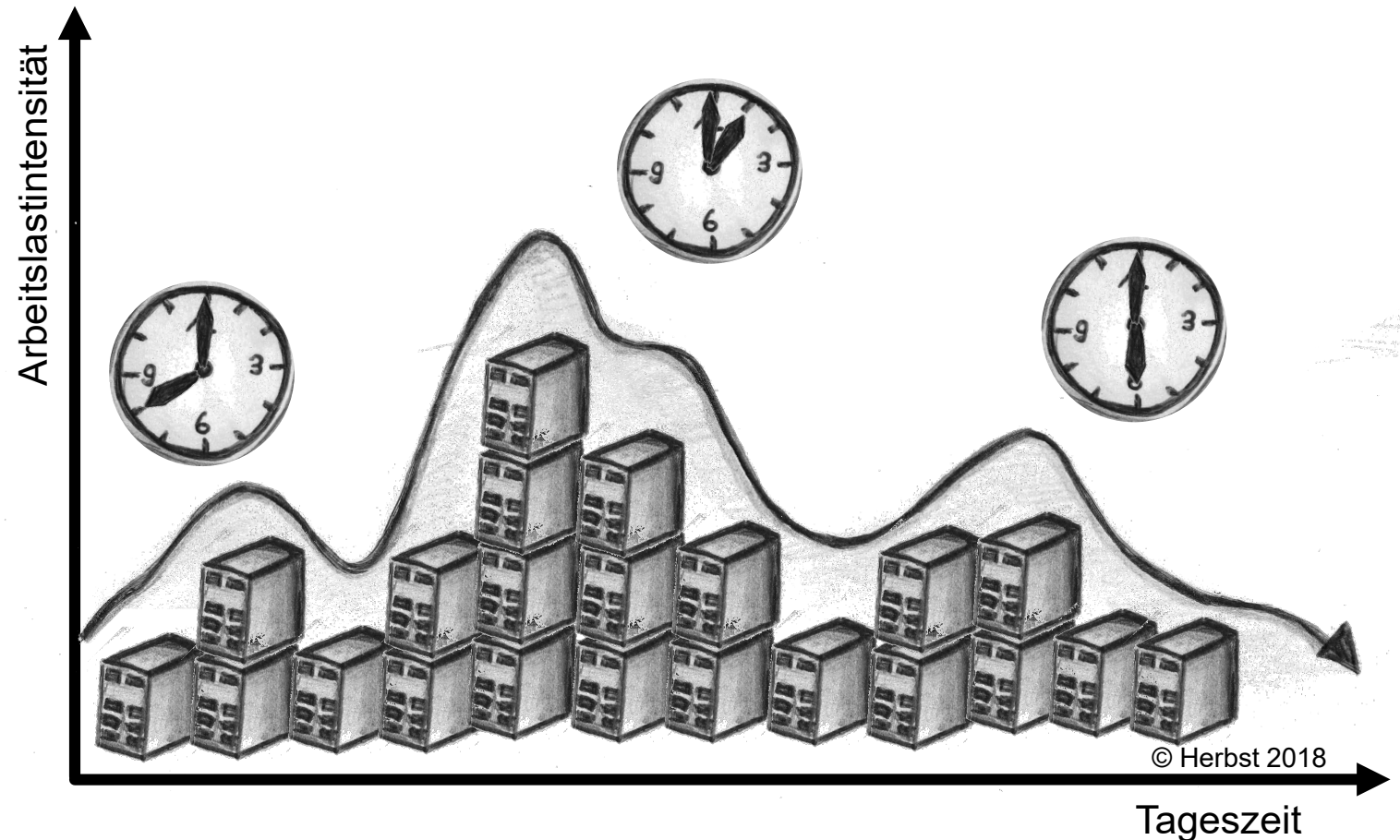
Prof. Dr. Samuel Kounev

Schloss Dagstuhl – 28. Mai, 2019

Automatisches
Hinzufügen und Entfernen
von virtuellen
Rechenressourcen

z.B.: VMs, vCPUs

Anpassen an wechselnde
Nachfrage über die Zeit
in der Größenordnung von
Minuten, Stunden oder Tagen

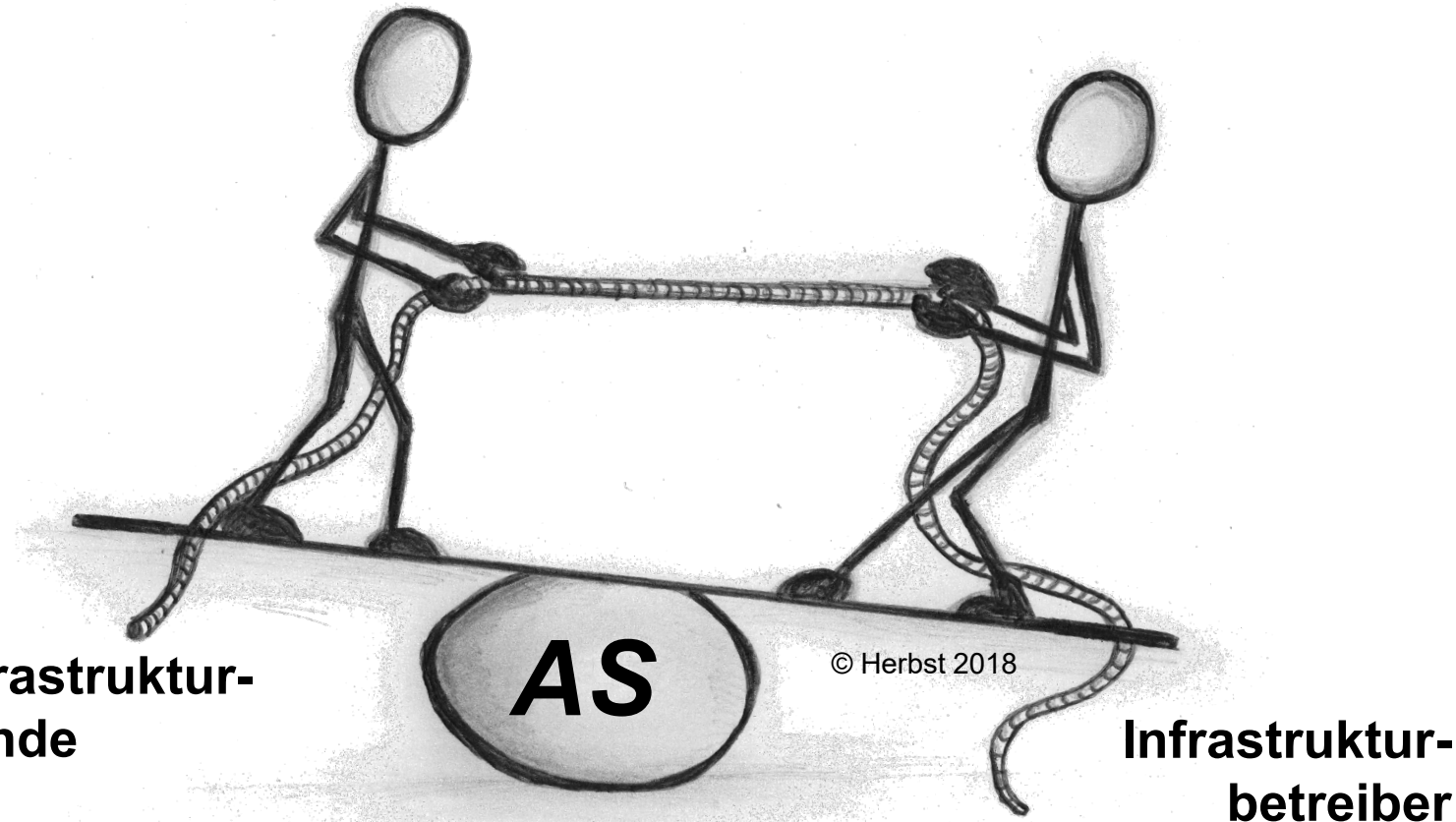


Automatisches Skalieren (AS) von Infrastruktur-Clouds

Automatisches Hinzufügen und Entfernen von virtuellen Rechenressourcen
z.B.: VMs, vCPUs

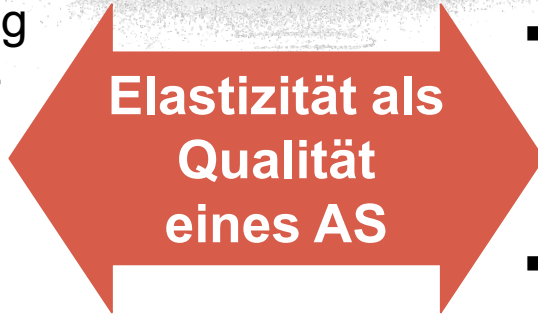
Anpassen an wechselnde Nachfrage über die Zeit in der Größenordnung von Minuten, Stunden oder Tagen

AS als Mediator zwischen Infrastrukturbetreiber und den Kunden



Infrastrukturkunde

- Stabile Leistung
- Exakt Bedarfsgerechter Ressourcenverbrauch



© Herbst 2018

Infrastrukturbetreiber

- Effizienter Betrieb ≈ Verkauf virtueller Ressourcen bei geringerer Nutzung physikalischer
- Zusätzliche Verwaltungsaufwände minimieren



“Elastizität im Cloud Computing ist der Grad zu dem ein System in der Lage ist, sich an **Änderungen in der Arbeitslast autonom anzupassen**, so dass sich **Allokation** und **Nachfrage** an Ressourcen zu **jedem Zeitpunkt** möglichst **optimal decken**.”

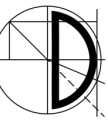
(sinngemäß übersetzt)

N. Herbst, S. Kounev and R. Reussner
Elasticity in Cloud Computing: What it is, and What it is
in Proceedings of the 10th International Conference on Auto-
Computing (ICAC 2013), San Jose, CA, June 24-28, 2013.



© Herbst 2018





Problem A



Kein Verfahren zur Elastizitätsbewertung

in realistischen Szenarien. Das führt zu

- a) hohem, **manuellem Aufwand**,
um einen AS optimal zu konfigurieren
- b) und **eingeschränkter Vergleichbarkeit**
sowie niedrigen **Wettbewerb** zwischen
AS Ansätzen

Ziel A

Definition eines **Benchmarks für AS**,
um so indirekt das Vertrauen in
neuartige, proaktive Ansätze und einen
breiteren Einsatz im Betrieb zu fördern.

Problem B



Hohes Risiko bei der Verwendung von AS

im Betrieb, da existierende Ansätze
entweder

- a) **rein reaktiv** arbeiten, ohne Möglichkeit
die nächsten Schritte zu planen, oder
- b) sich verlassen auf **einzelne, prädiktive
Modelle** angewandt **in Isolation**.

Ziel B

Reduzierung des Risikos einen neu-
artigen AS im Betrieb zu nutzen
durch den Einsatz **mehrerer proaktiver
Mechanismen** in Kombination mit
einem reaktivem Rückfall-Mechanismus.

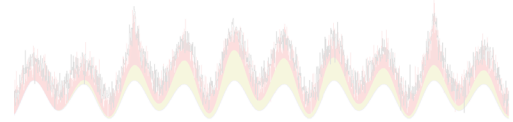


Ziel A

Definition eines **Benchmarks für AS**, um so indirekt das Vertrauen in neuartige, proaktive Ansätze und einen breiteren Einsatz im Betrieb zu fördern.

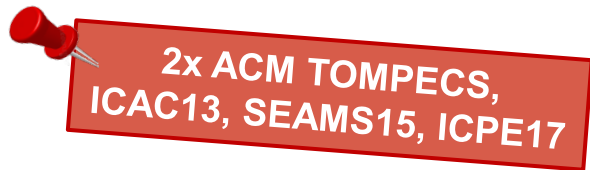
1

LIMBO: Werkzeuggestützte od. automatische **Modellierung** von Lastintensitätsprofilen



2

BUNGEE Messmethodik und Metriken um **Elastizität** in **fairer und wiederholbarer** Weise zu bewerten.



Ziel B

Reduzierung des Risikos einen neuartigen AS im Betrieb zu nutzen durch den Einsatz **mehrerer proaktiver Mechanismen** in Kombination mit einem reaktivem Rückfall-Mechanismus.

3

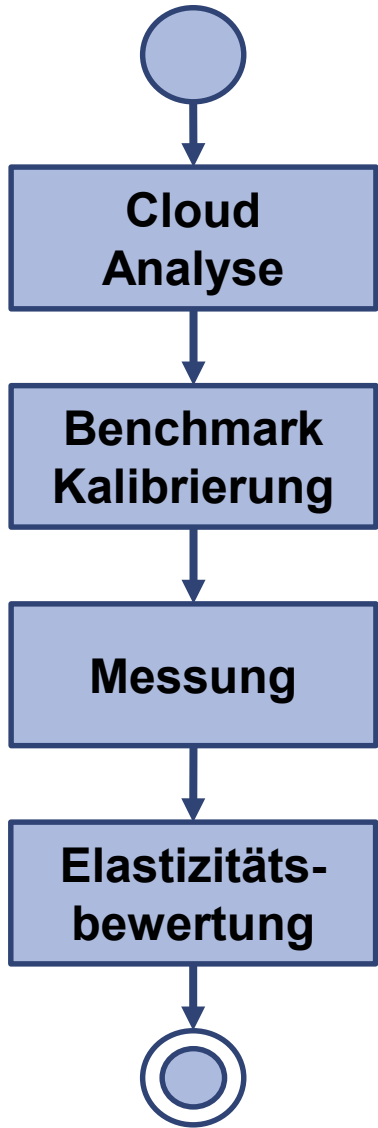
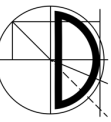
Chameleon – hybrider Auto-Skalierer, der geschickt proaktive und reaktive Methoden vereint

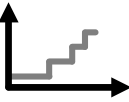


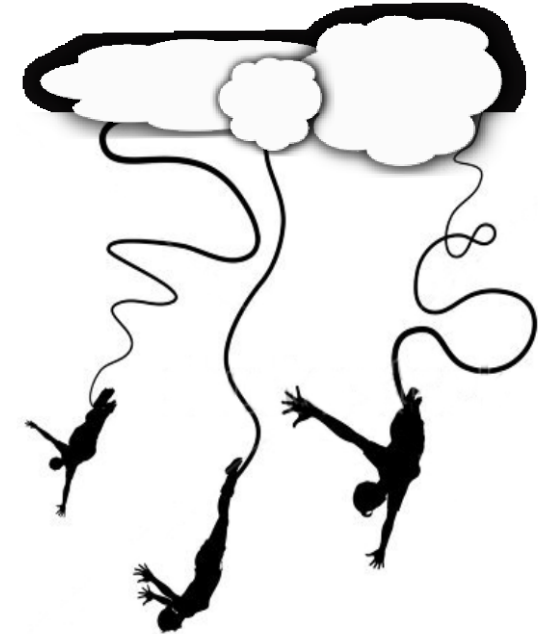
4

Telescope – Zerlegungs-basierte Vorhersage von Zeitreihen, um zuverlässiger genaue Vorhersagen pünktlich zur Verfügung zu stellen

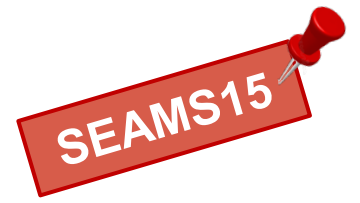




Performanzanalyse der zugrunde liegenden Rechenressourcen und deren Skalierbarkeit 



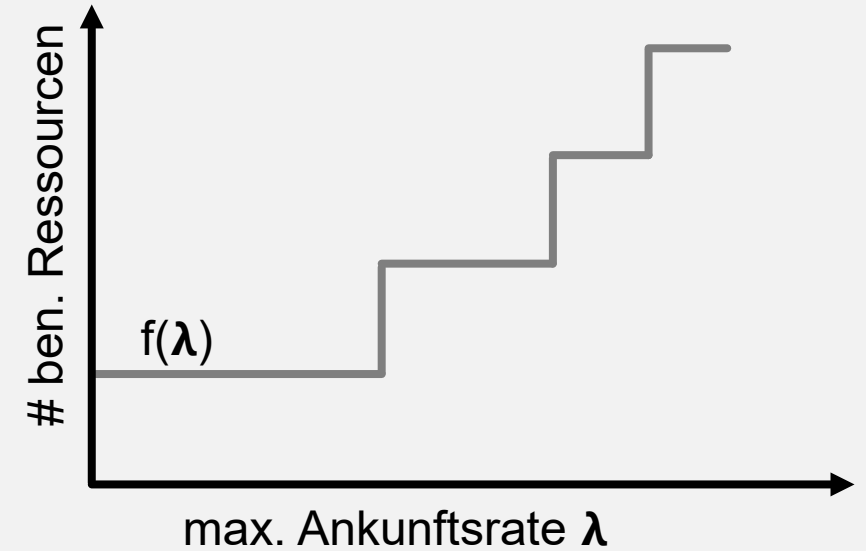
Cloud Elasticity Benchmark
BUNGEE



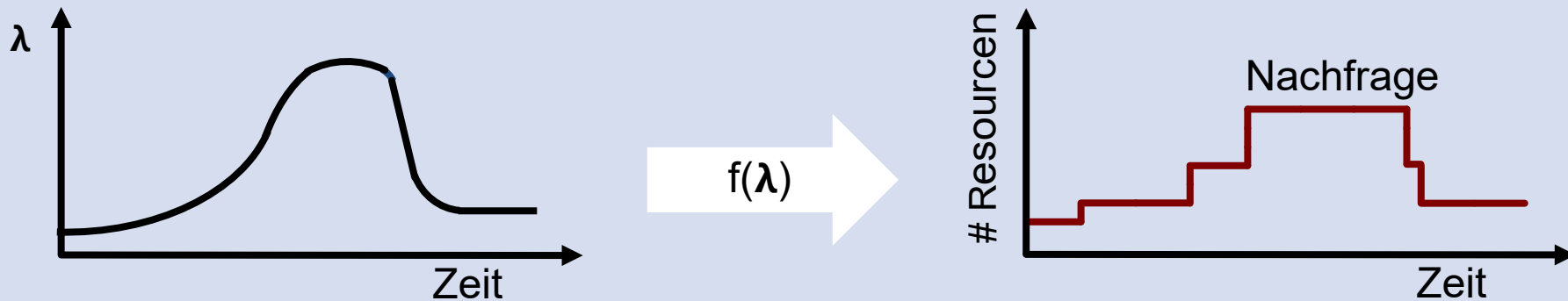


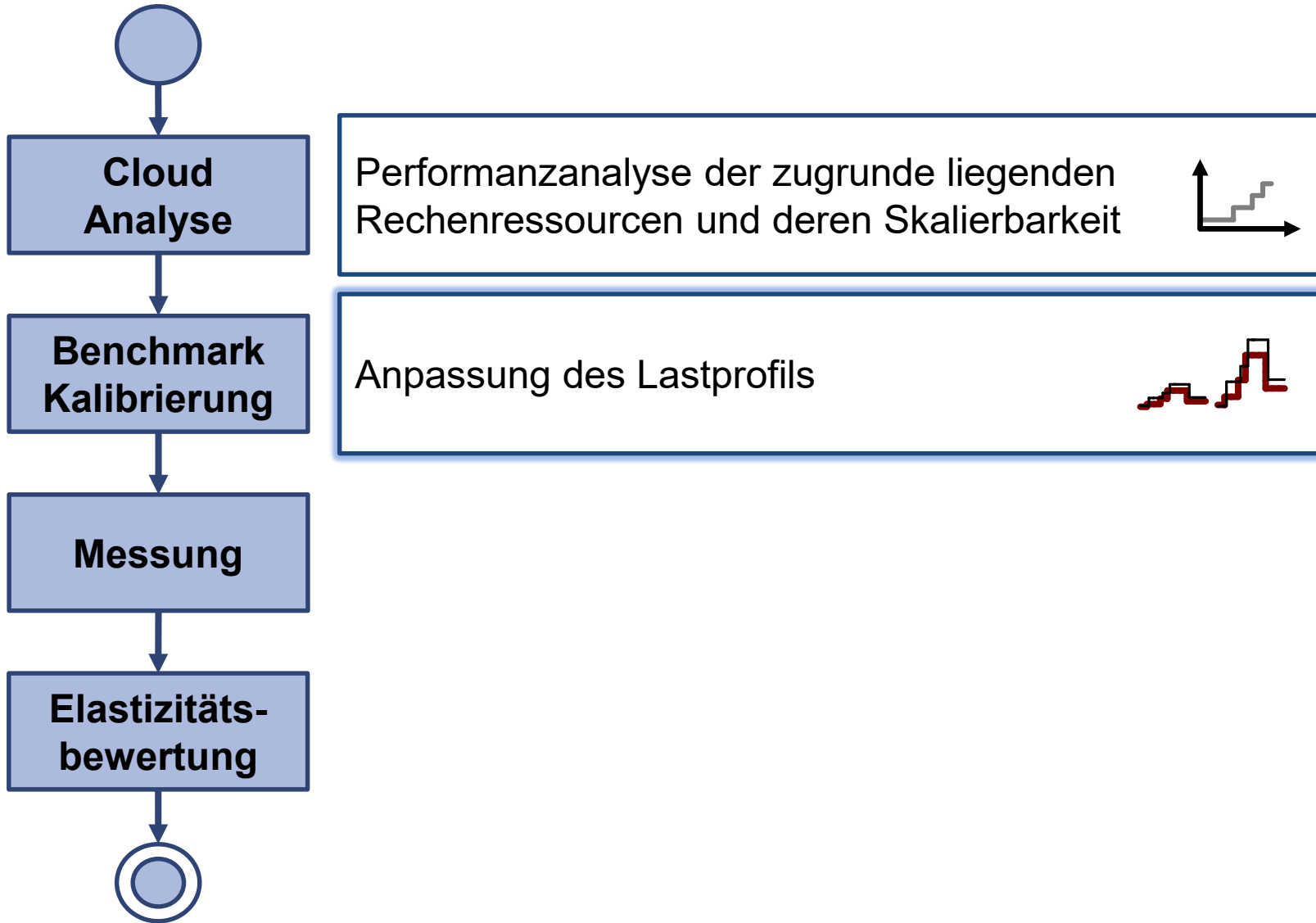
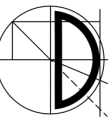
Grundidee:

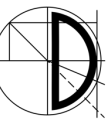
- (1) Bewerte die Leistungsfähigkeit jeder Stufe
- (2) Finde Kapazität = **max. Ankunftsrate λ** ,
der die Cloud Anwendung standhält
ohne Dienstgütereletzungen,
dazu binäre Suche
- (3) Stufenfunktion für die Nachfrage:
 $f(\lambda) = \#$ benötigter Ressourceneinheiten



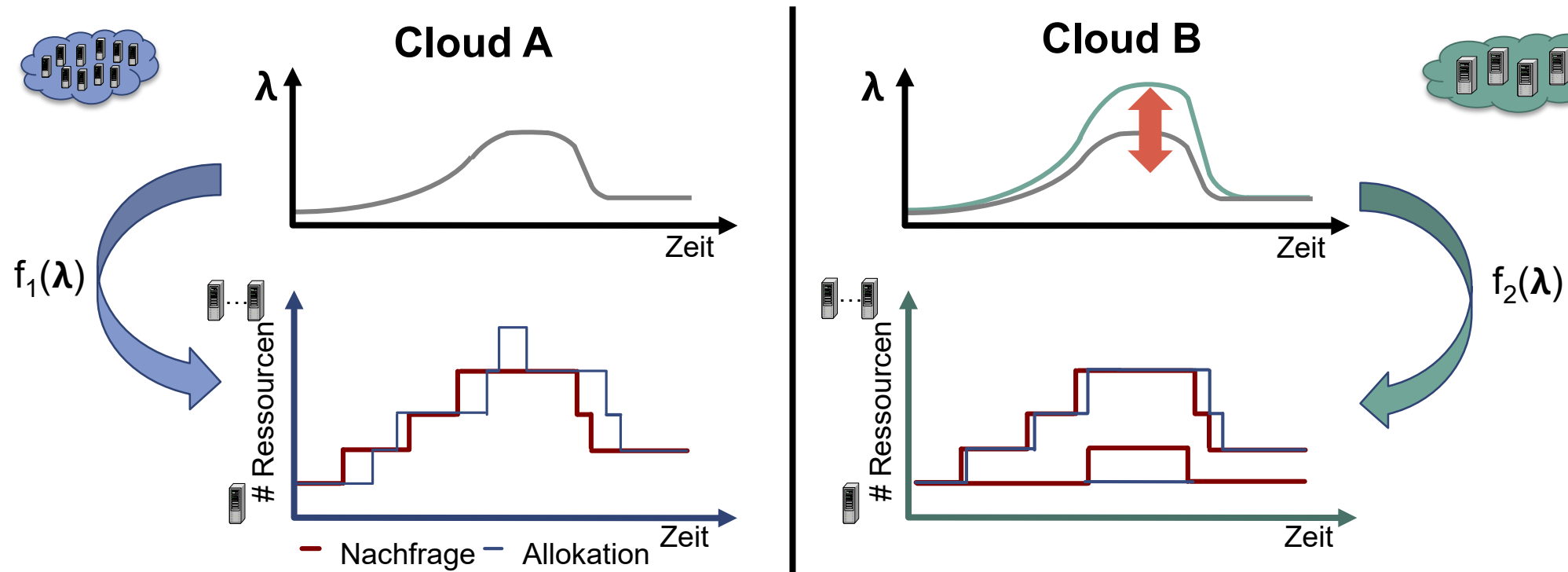
Ableiten der tatsächlichen Ressourcennachfrage für beliebige Lastintensitätsprofile als optimaler Auto-Skalierer (Goldstandard)





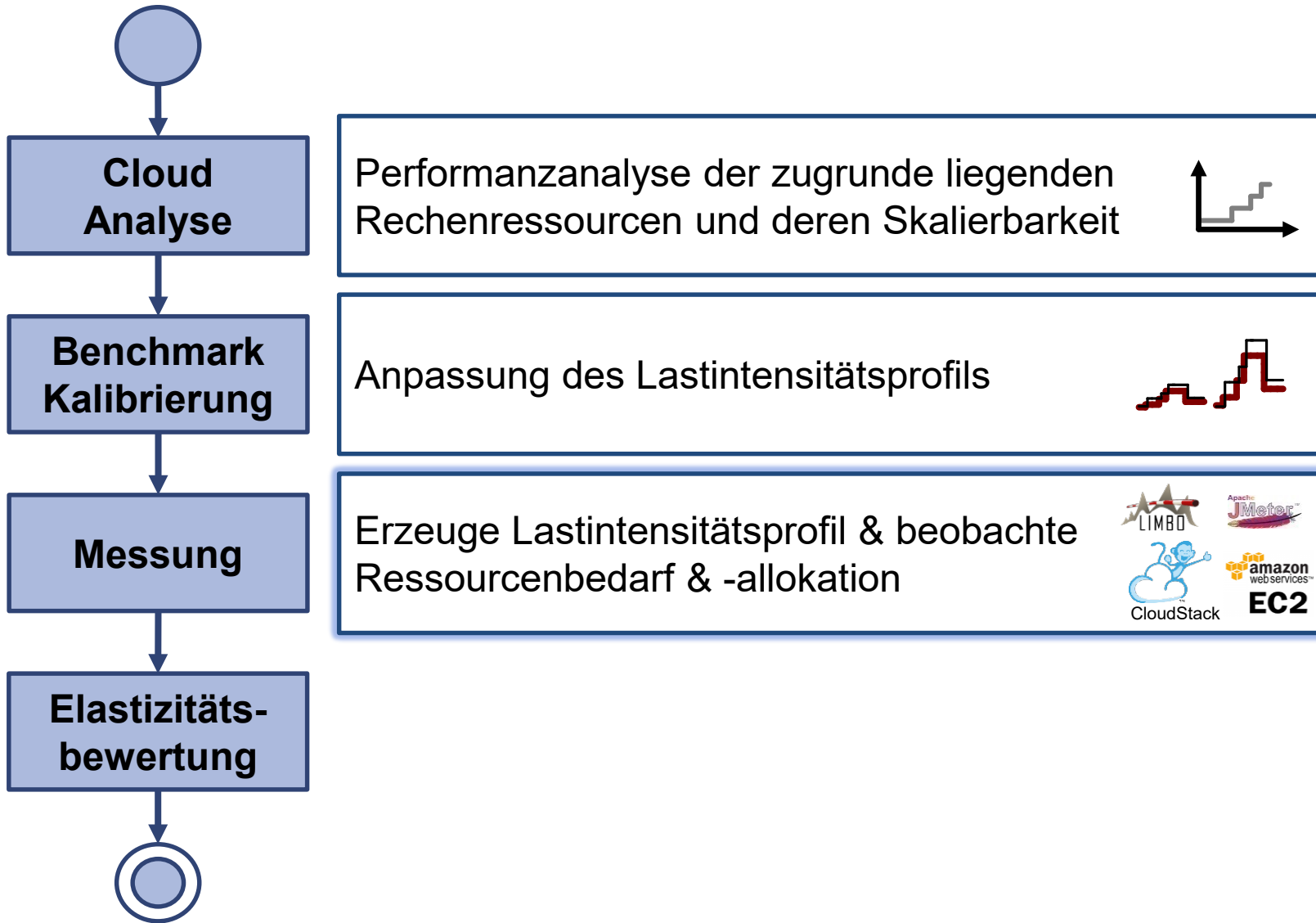


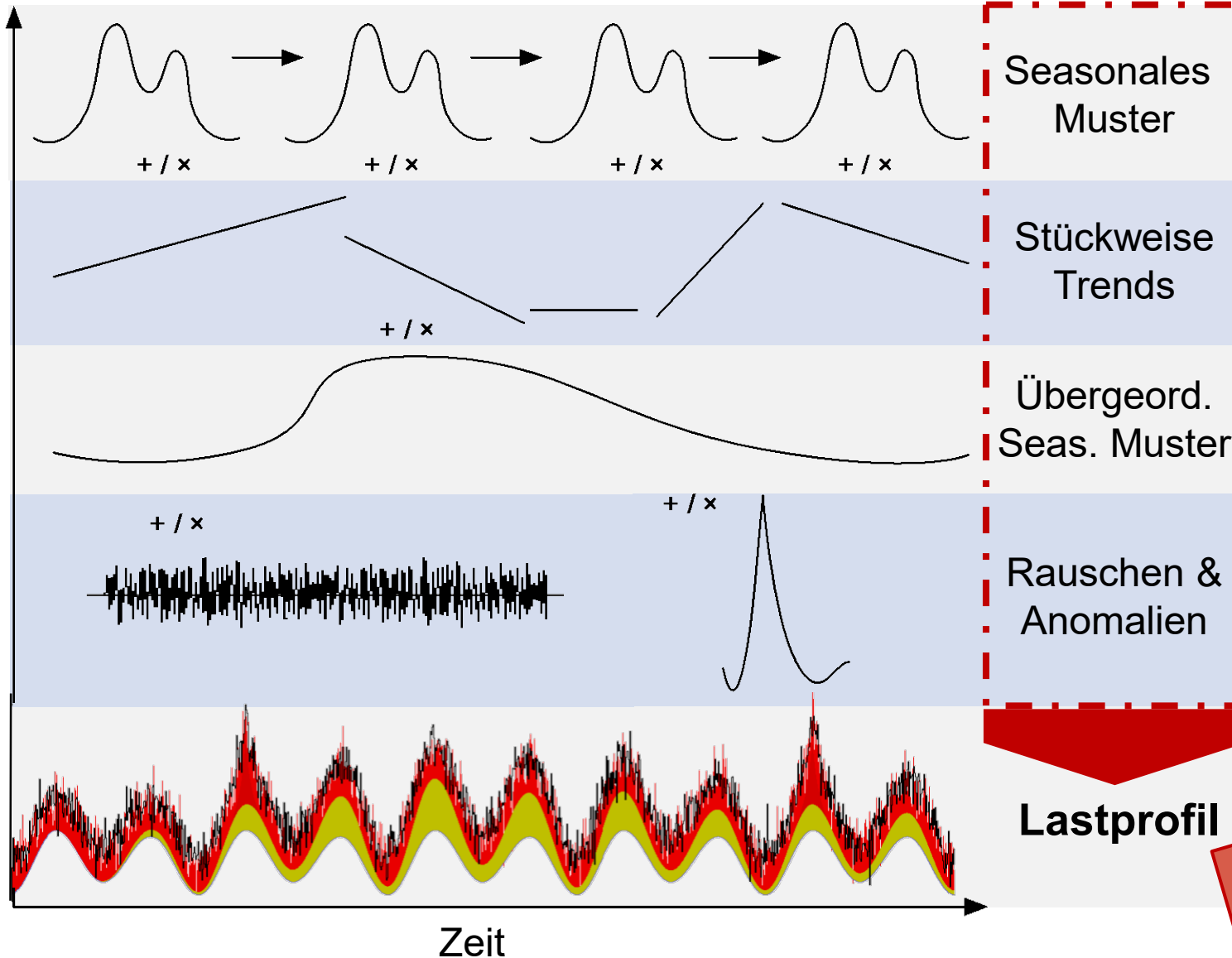
Beabsichtigt: Erzeuge dieselbe Nachfrage an Ressourcen in untersch. Clouds



Lösungsweg: Anpassung des Lastintensitätsprofils unter Berücksichtigung

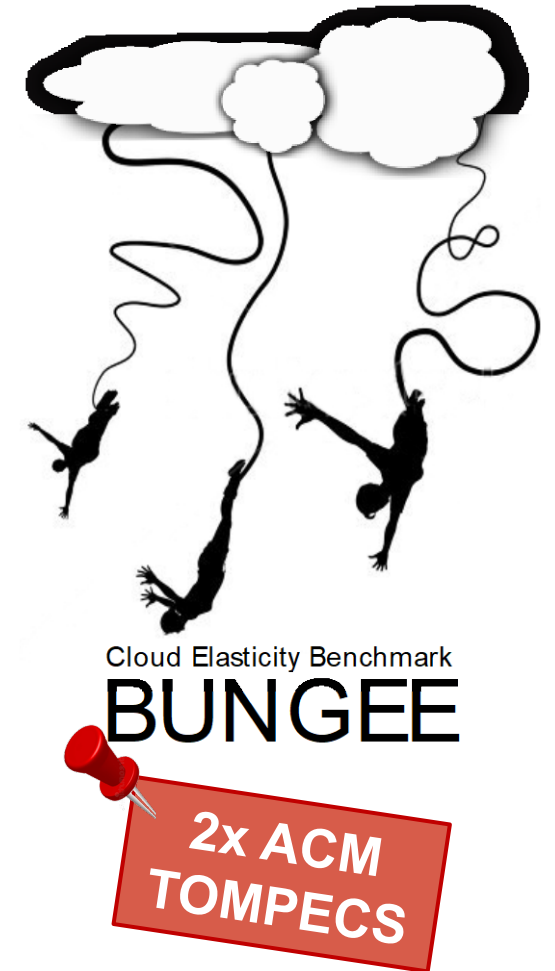
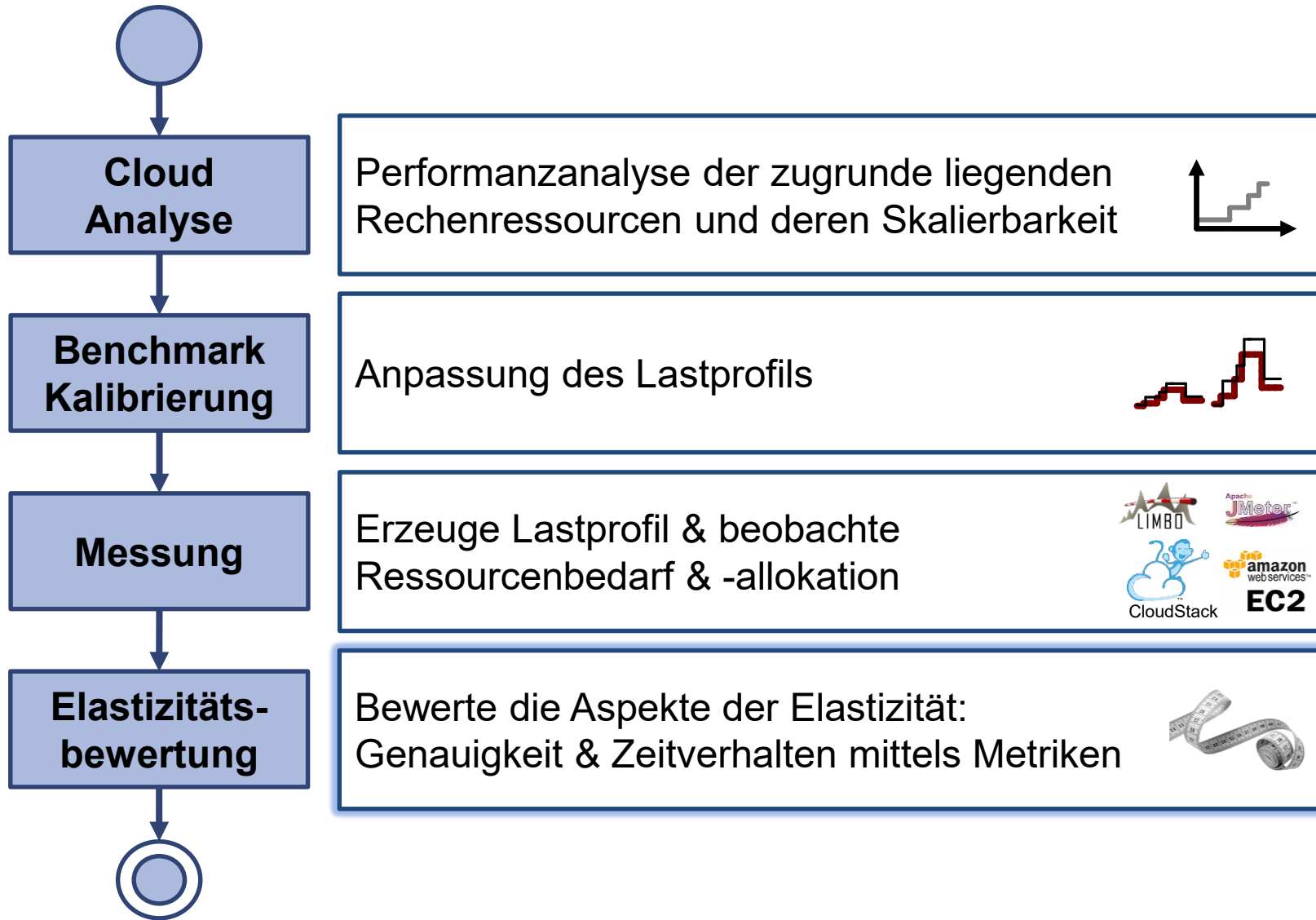
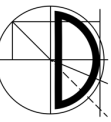
- der unterschiedlichen Leistungsfähigkeit der Ressourcen und
- unterschiedlichem Skalierungsverlauf

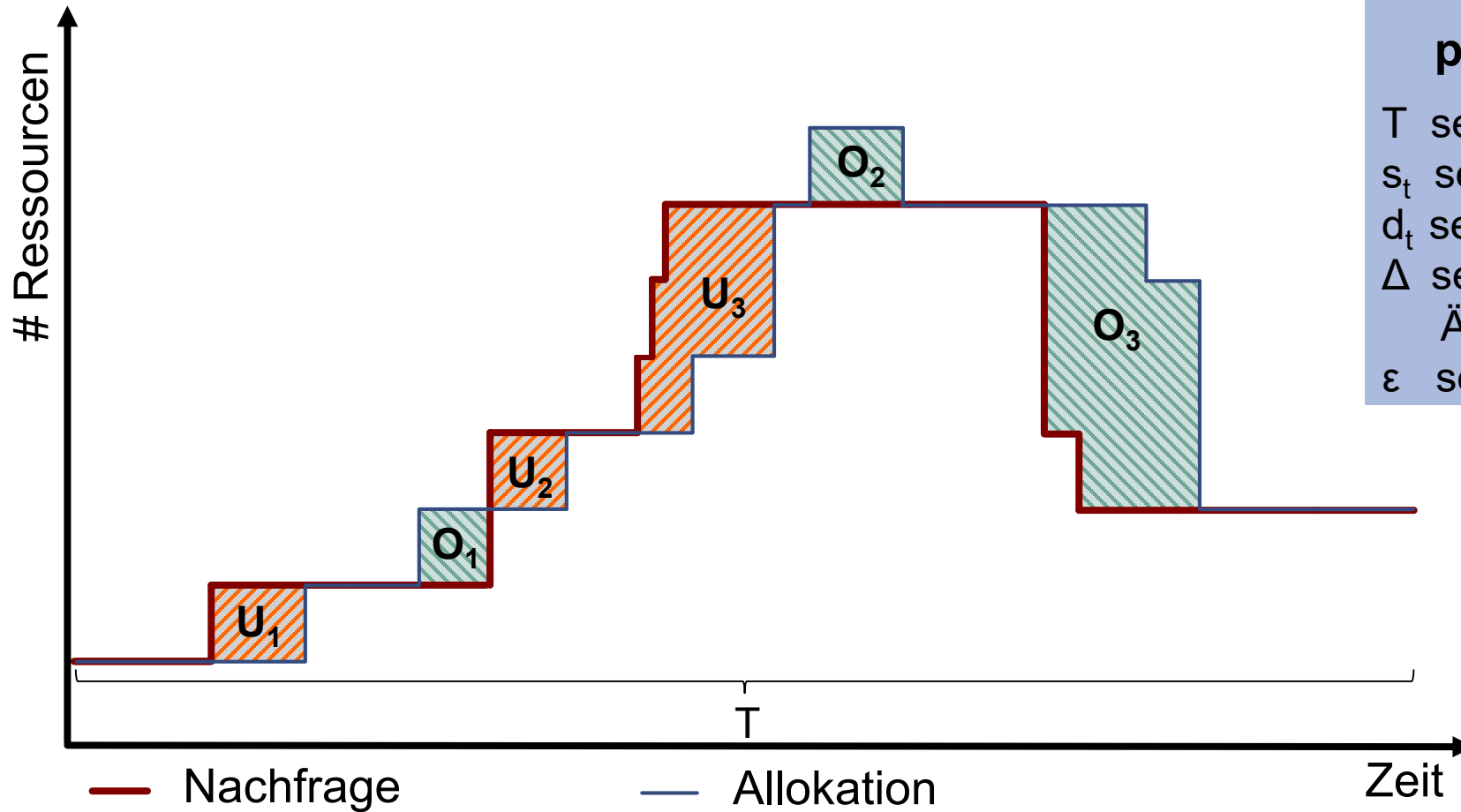




ACM TAAS 2017,
SPEC RG peer-reviewed tool







∅ Anteil θ / Menge a an unter- bzw. überprovisionierten Ressourcen

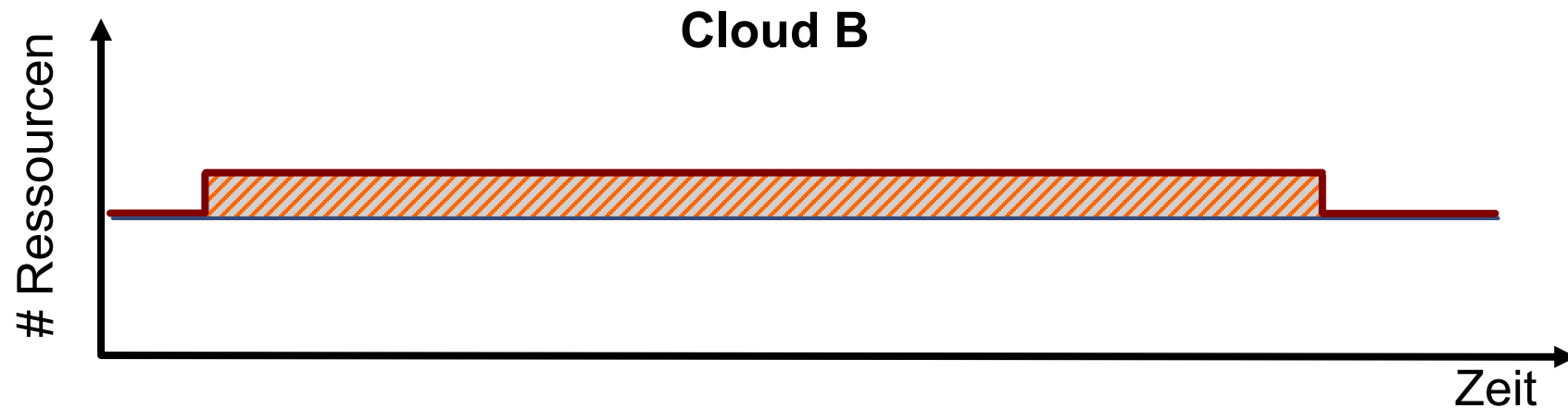
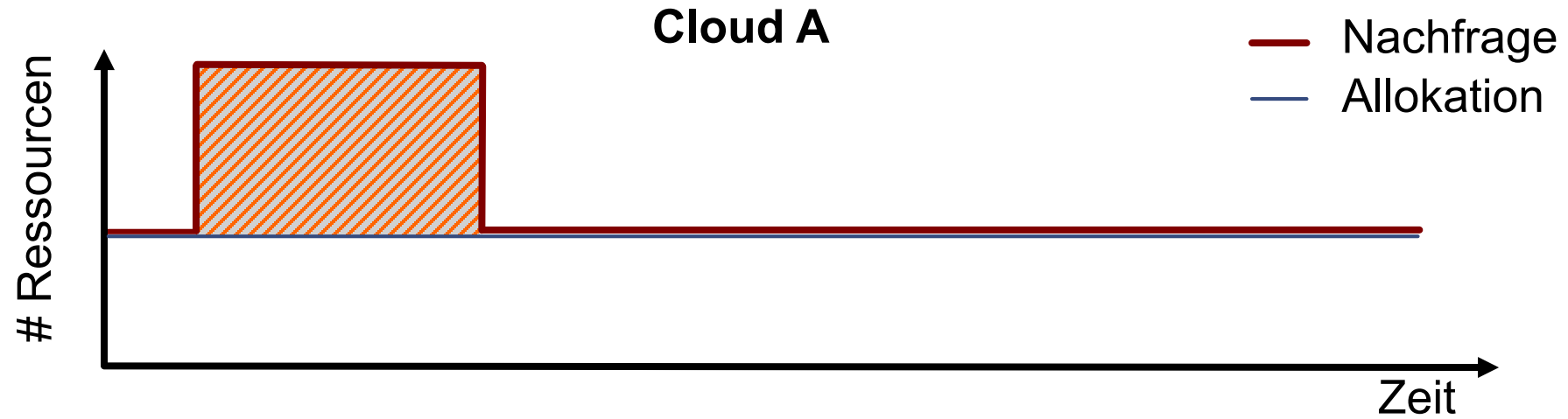
T sei die Experimentdauer
 s_t sei die Allokation zum Zeitpunkt t
 d_t sei die Nachfrage zum Zeitpunkt t
 Δ sei Zeit zw. letzter und aktueller Änderung
 ϵ sei min. Nachfrage, z.B. $\epsilon = 1$

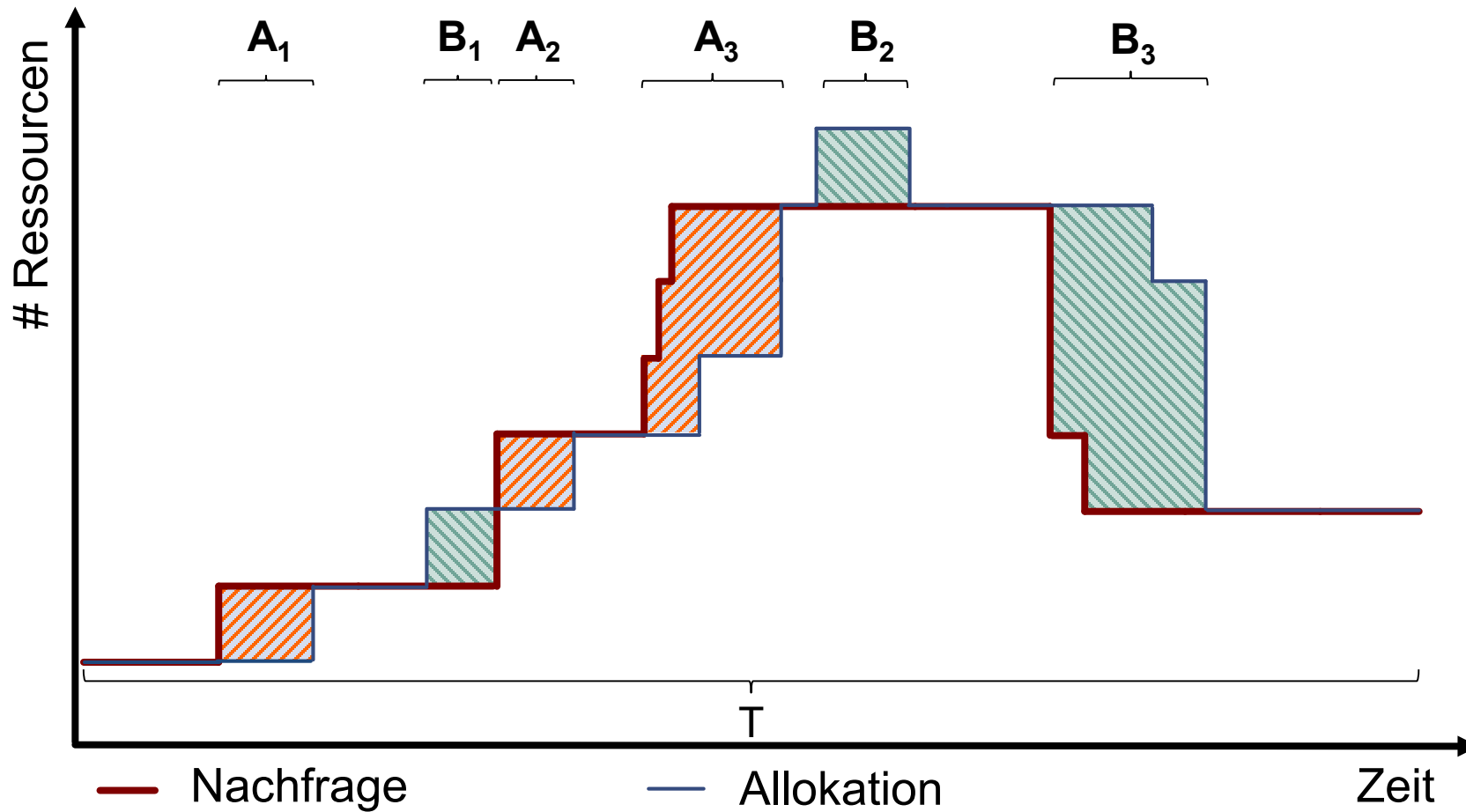
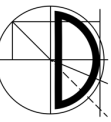
$$\theta_U[\%] := \frac{100}{T} \cdot \sum_{t=1}^T \frac{\max(d_t - s_t, 0)}{\max(d_t, \epsilon)} \Delta t$$

$$\theta_O[\%] := \frac{100}{T} \cdot \sum_{t=1}^T \frac{\max(s_t - d_t, 0)}{\max(d_t, \epsilon)} \Delta t$$

$$a_U[\#res] := \frac{1}{T} \cdot \sum_{t=1}^T \max(d_t - s_t, 0) \Delta t$$

$$a_O[\#res] := \frac{1}{T} \cdot \sum_{t=1}^T \max(s_t - d_t, 0) \Delta t$$



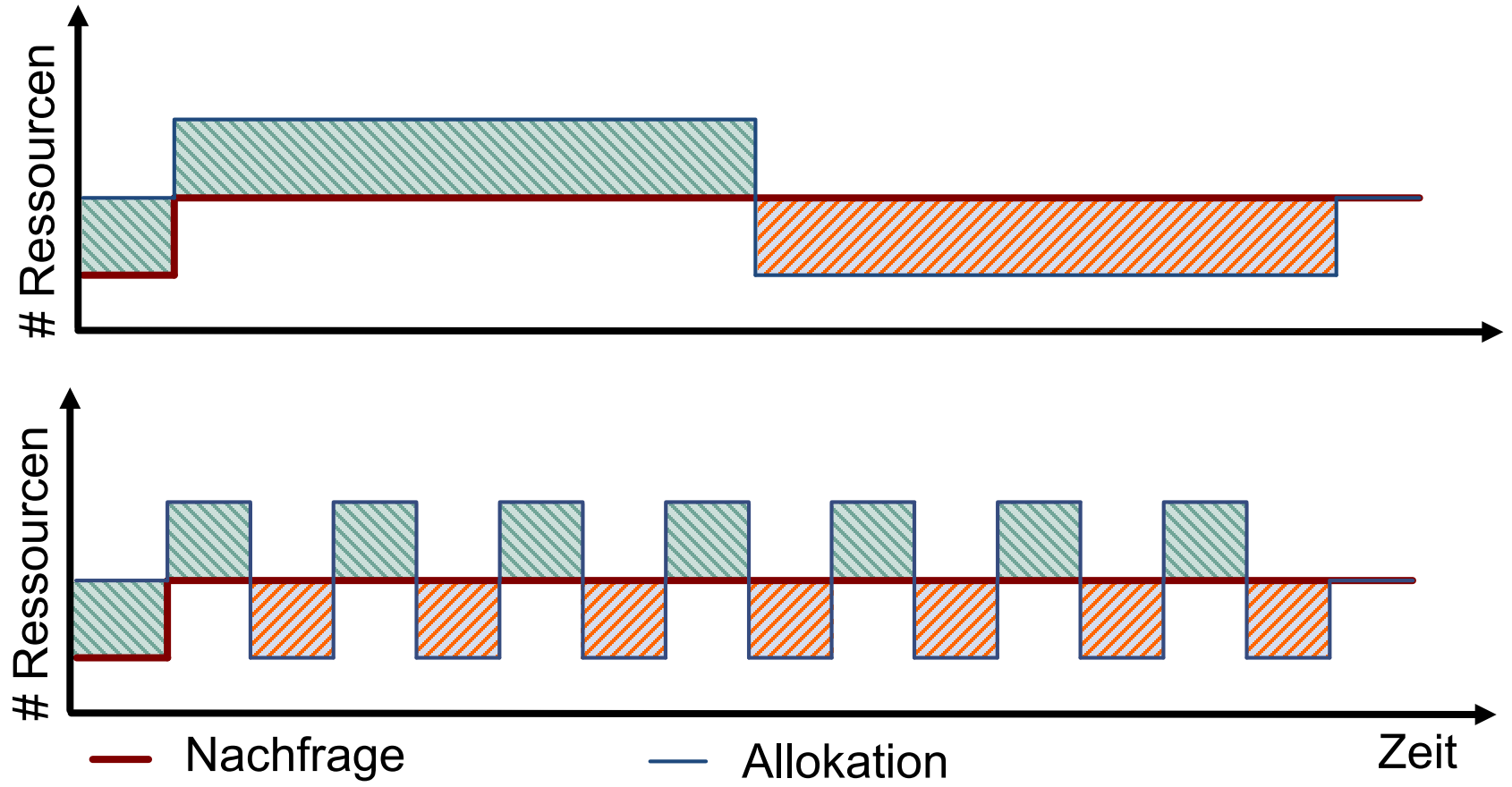


∅ Zeit τ in unter- bzw. überprovisioniertem Zustand

T sei die Experimentdauer
 s_t sei die Allokation zur Zeit t
 d_t sei die Nachfrage zur Zeit t
 Δ sei Zeit zw. letzter und aktueller Änderung

$$\tau_U[\%] := \frac{100}{T} \cdot \sum_{t=1}^T \max(\text{sgn}(d_t - s_t), 0) \Delta t$$

$$\tau_O[\%] := \frac{100}{T} \cdot \sum_{t=1}^T \max(\text{sgn}(s_t - d_t), 0) \Delta t$$



Anteil der Zeit u in der Allokation und Nachfrage nicht parallel verlaufen

T sei die Experimentdauer
 s_t sei die Allokation zur Zeit t
 d_t sei die Nachfrage zur Zeit t
 Δ sei Zeit zw. letztem und aktueller Änderung

$), 1) \Delta t$

→ Intuitive Interpretation durch Wertebereich $[0, 1]$

(1) Verwendung von LIMBO Lastprofilmodellen

1

- Extrahiert von 5 verschiedenen realen Aufzeichnungen
FIFA WC98, BibSonomy, IBM CICS, German Wiki, Retraire
- Generierung von 1-3 Tagen jeweils pro Experimentzeit pro Tag
- Eine Matrixzerlegung mit der RT Suite

(2) Anwendung der

2

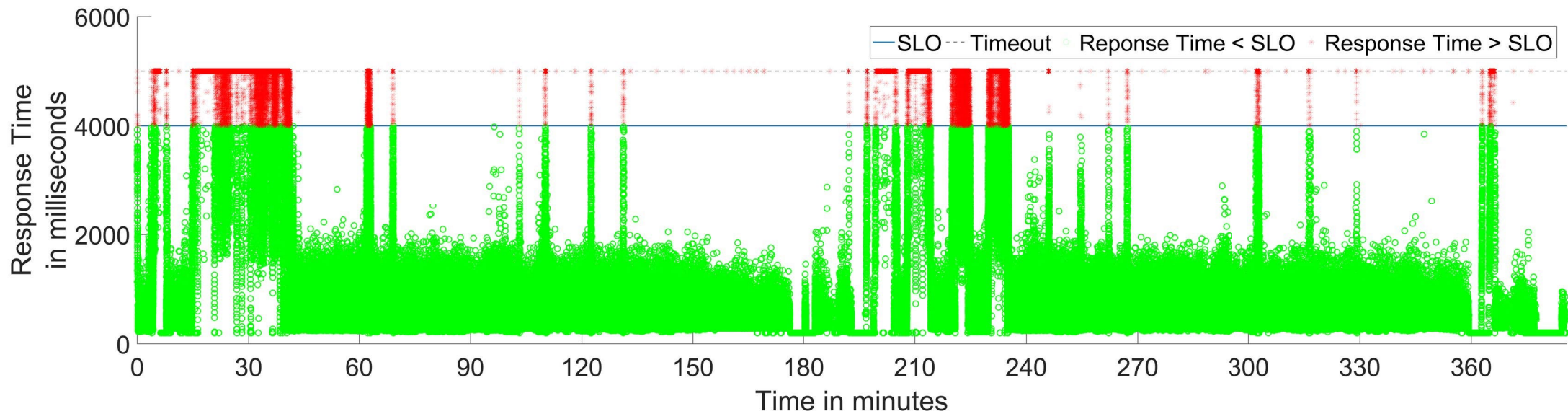
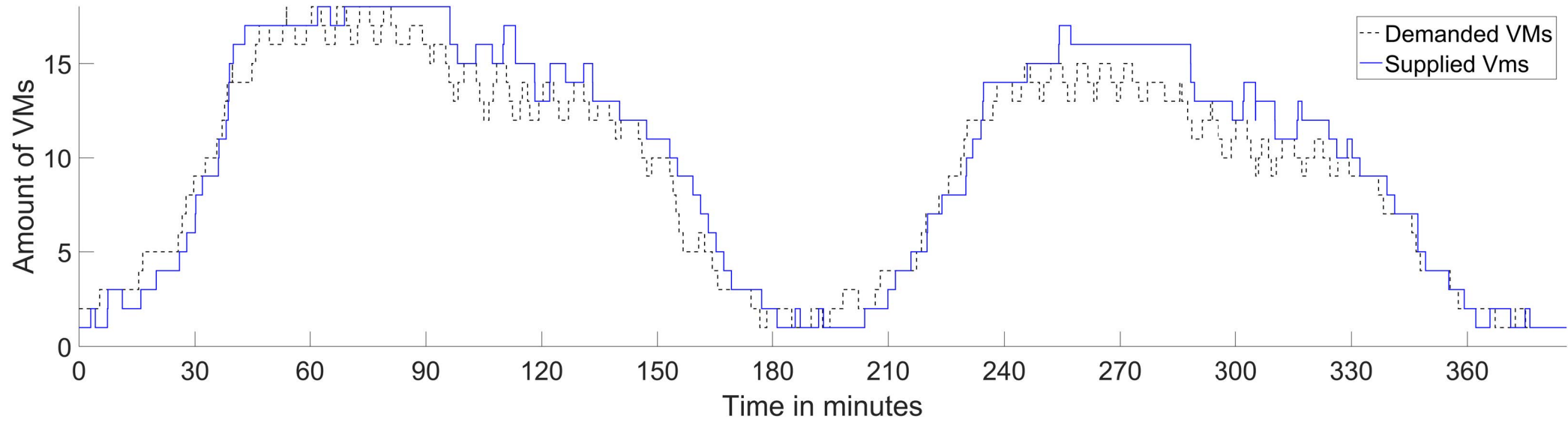
- Auf Installationen in
privates Apache C
- Mit einer oberen Gr
ds:
e Forschungscloud)
2-8GB Arbeitsspeicher

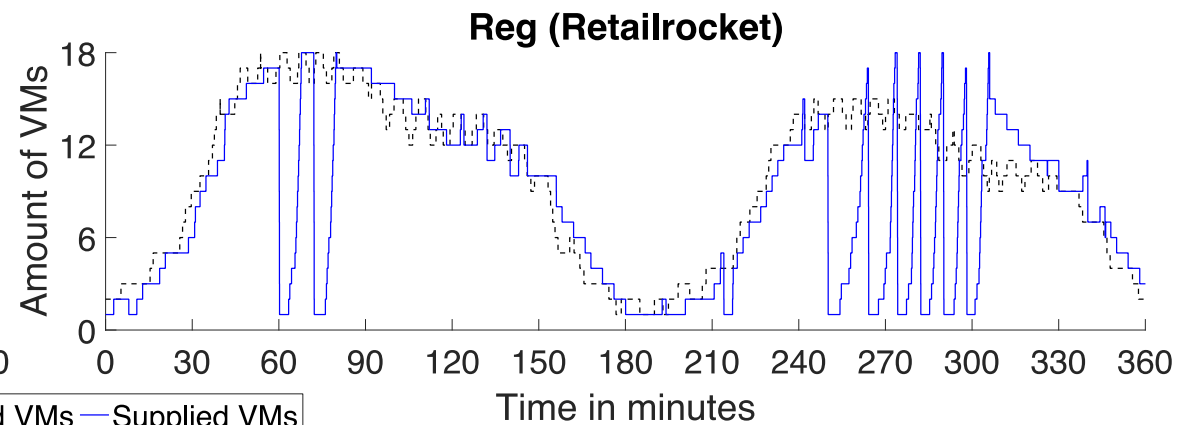
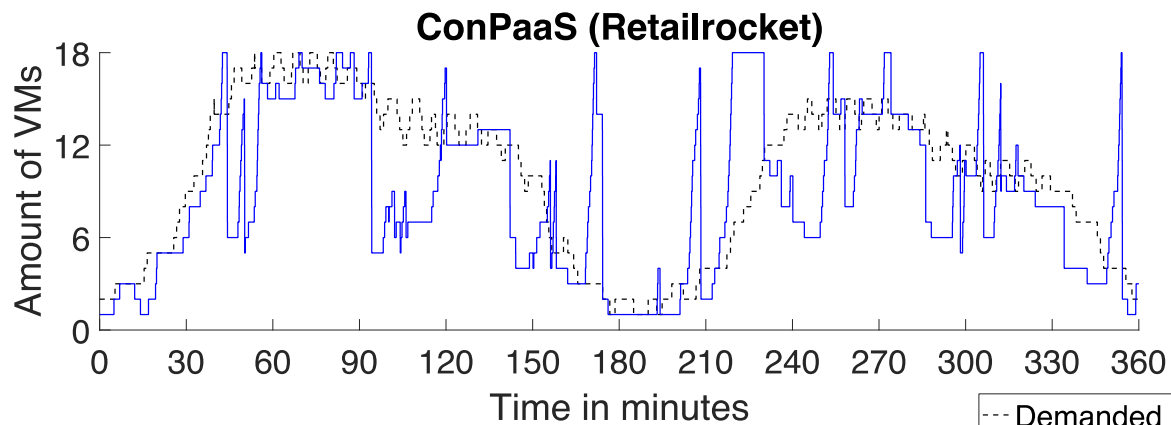
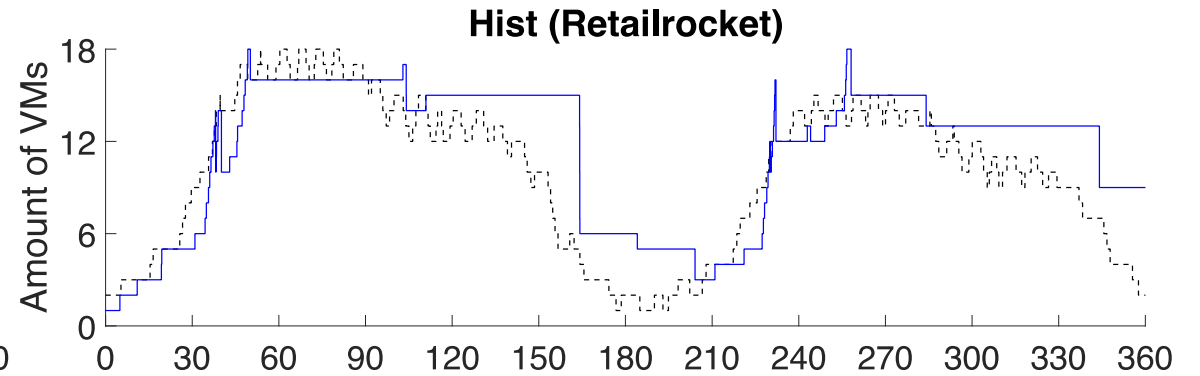
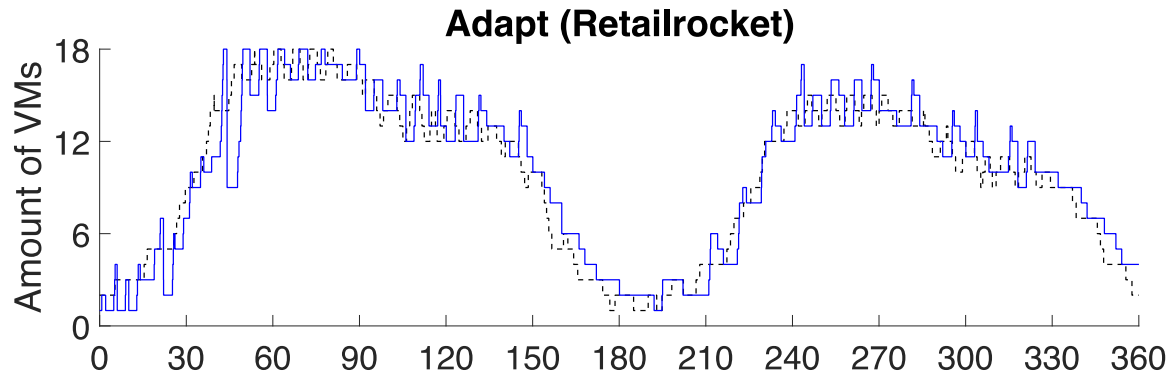
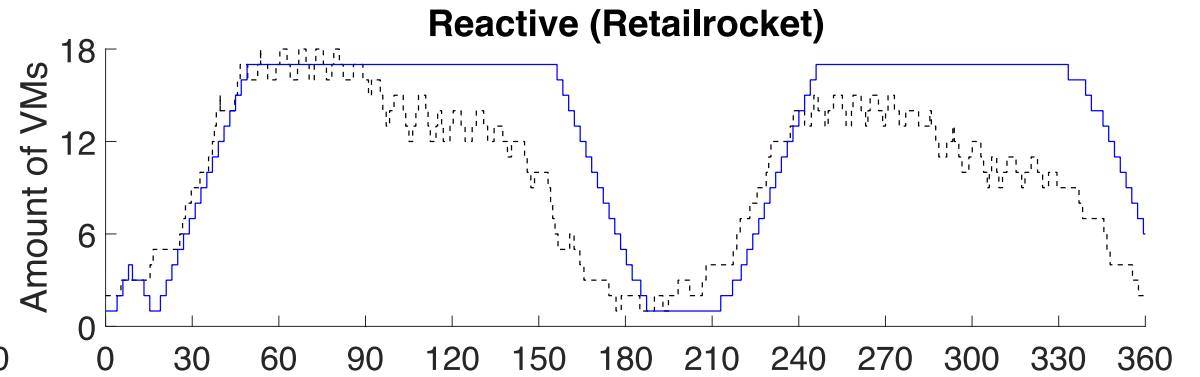
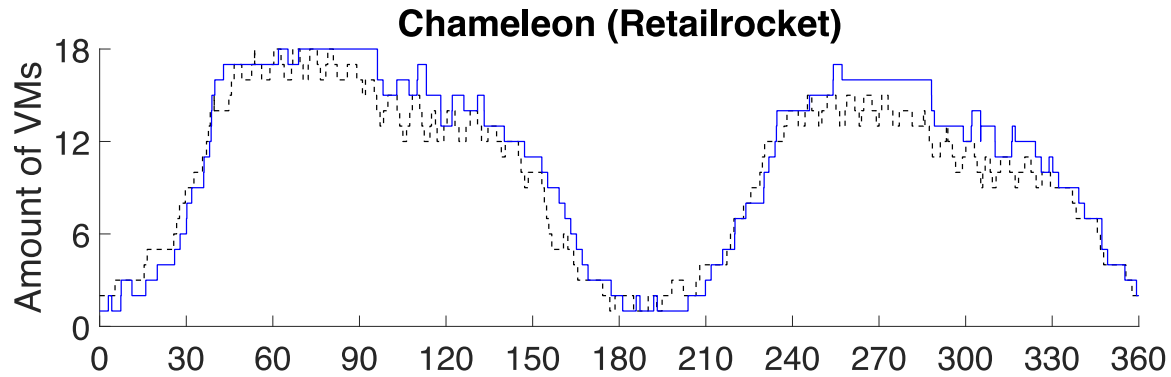
(3) Wettbewerb des

3

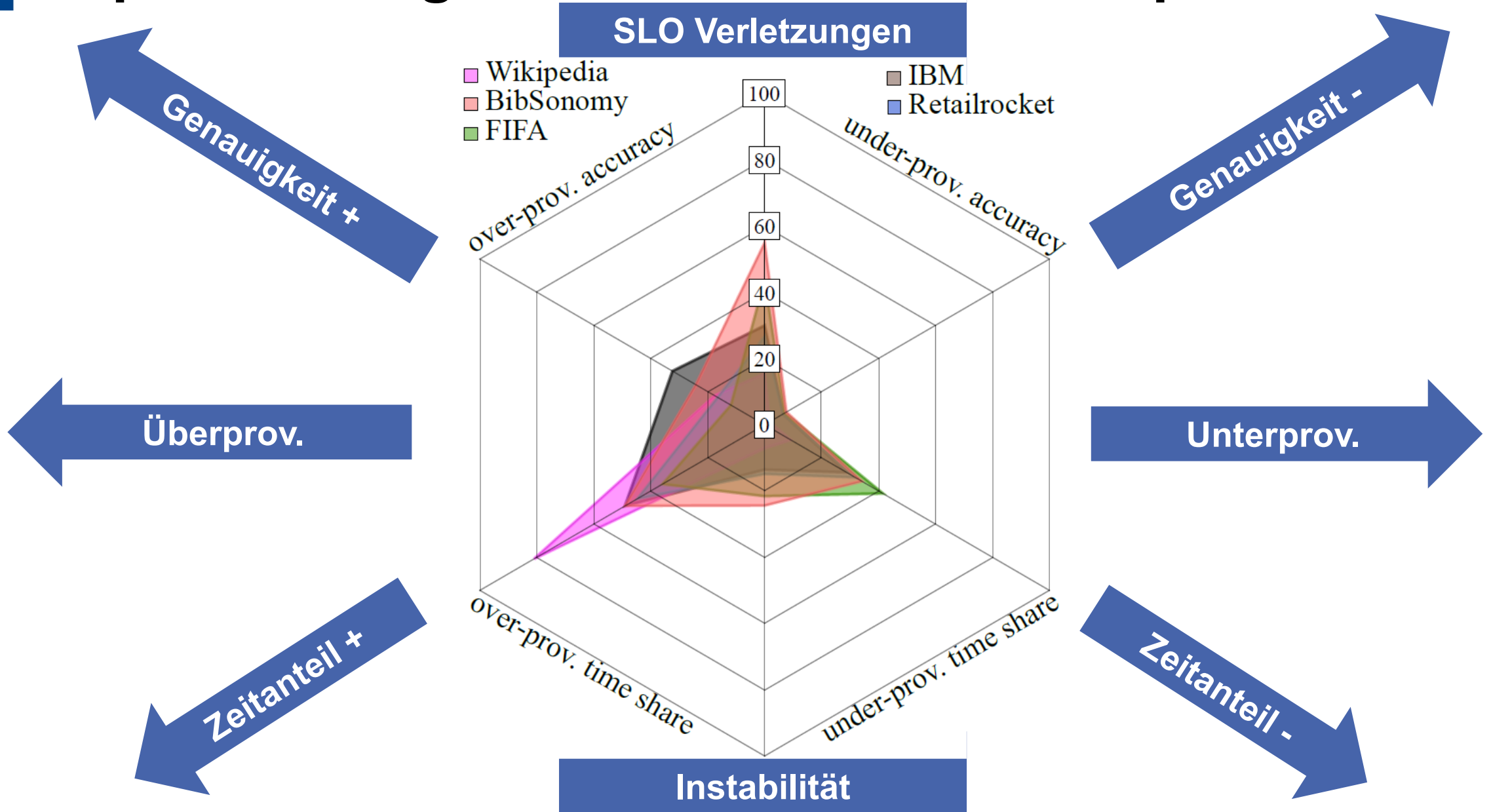
- gegen 4 proactive A
- Adapt [Adhikari12] aus
- Hist [Urgaonkar08] aus
- Reg [Iqbal11] und ConP
der Zeitreihenanalyse, sowie
- 1 reaktiver Auto-Skalierer, sowie einem Szenario mit
statischer Ressourcenallokation.

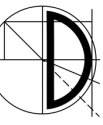
Evaluationsumfang
 400 Stunden an Experimenten
 107 Mio. generierte Anfragen
 5000 Anpassungen
 ~700 € bei AWS
 pro Experimentstunde:
 ~0.27 Mio. Anfragen, und
 13 Anpassungen im Durchsch.



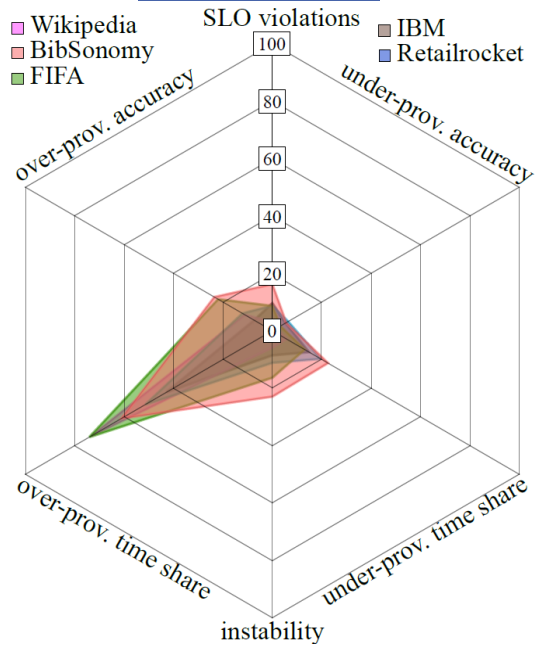


--- Demanded VMs — Supplied VMs

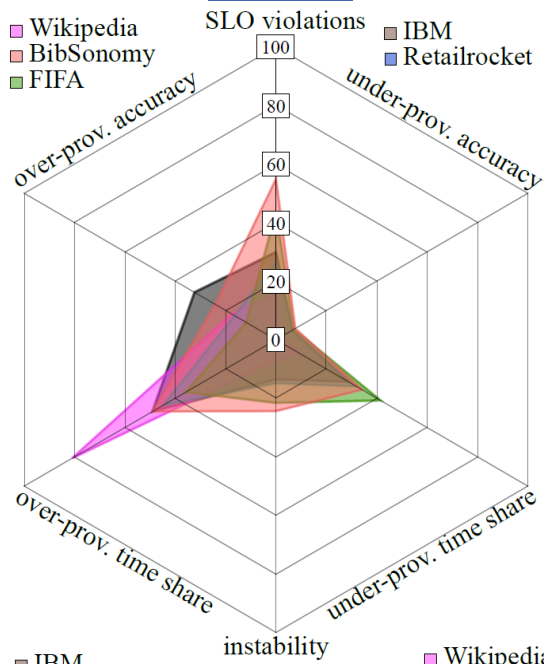




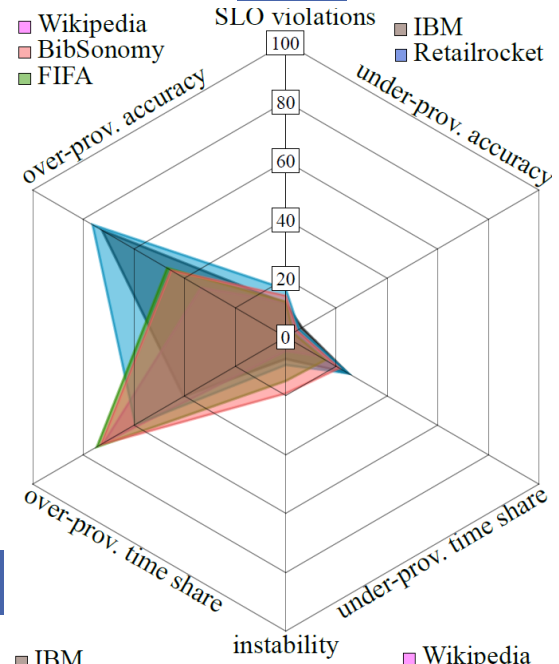
Chameleon



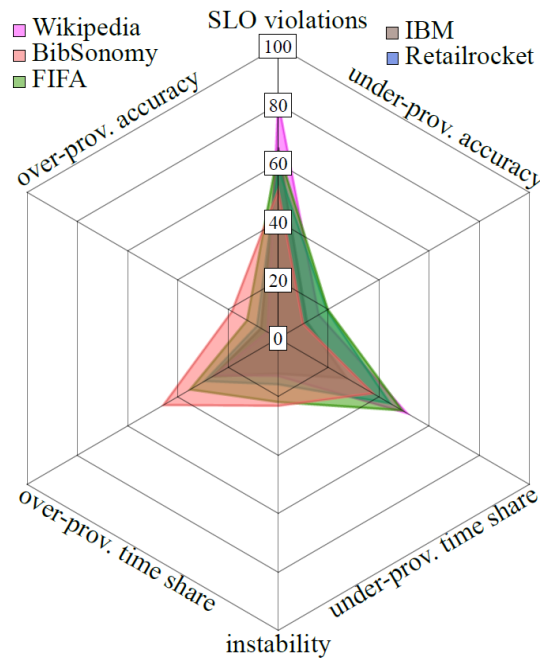
Adapt



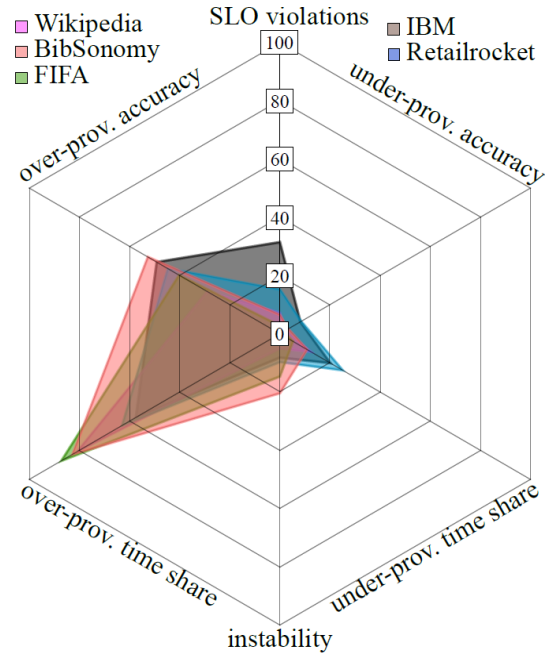
Hist



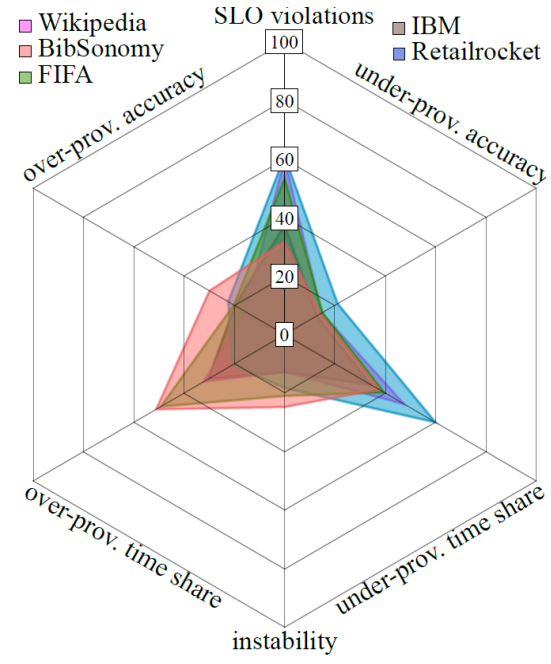
Reg



Reactive



ConPaaS



Erreichen von Ziel A

Definition eines **Benchmark** für Auto-Skal.

- **LIMBO Lastintensitätsprofile**
- **BUNGEE Elastizitätsmessmethodik und Metriken**

1

2

Erreichen von Ziel B

Risikoreduzierung beim Einsatz neuartiger Auto-Skalierer

- **Chameleon: Hybrider Auto-Skalierer** geht als **Sieger** eines breit angelegten **Wettbewerbs** hervor
- **Telescope: Zerlegungsbasierte Vorhersagemethode** verbessert AS

3

4

Gestartete und zukünftige Forschung

“Herausforderungen des Ressourcenmanagements durch neuste Cloud-Trends”

- Containerisierung, Micro-Dienste, Function-as-a-Service & Serverless, Fog/Edge
 - Ressourcenverschachtelung und Fragmentierung
 - Ressourcenprioritäten, Reservierungen und kostenbewusstes Management
 - Neuauflage des Platzierungsproblems
- Anwendungsänderungen durch DevOps
- Verbesserte Lastvorhersage
- Benchmarkmethodik und Metriken für mehr-dimensionale Elastizität / Ressourcenmanagement



© Herbst 2018



Sichtbarkeit der Beiträge

- ✓ 5 Zeitschriftartikel Wiley CCPE, 2 ACM TOMPECS, ACM TAAS, IEEE TPDS
- ✓ 3 Artikel in IEEE TSE, IEEE Internet Comp. Revision IEEE Software
- ✓ 7 Konferenz- ACM/SPEC ICPE13,17,18 langbeiträge 2x SEAMS15, IEEE Cloud
- ✓ 4 Kurzbeiträge ICAC13,17, Closer18, ITISE17
- ✓ 2 Buchkapitel Springer SAC und mehr...

4 neue, quelloffene Software-Werkzeuge @ <https://descartes.tools>

- IBM PhD Fellowship
- "SPECTacular award for tec. leadership"
- Unibund Förderpreis der Mainfränkischen Wirtschaft

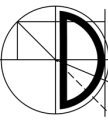


9 Masterarbeiten betreut,
5 davon Promovenden,
4 davon Preisträger

2x Workshop-organisator:

SeAC
Self-Aware Computing
@ FAS*

HotCloudPerf 2018
hotcloudperf.spec.org
Berlin, Germany, April 9
@ ACM/SPEC ICPE



Adhikari12 R. Adhikari and R. Agrawal, "An Introductory Study on Time Series Modeling and Forecasting," arXiv preprint arXiv:1302.6613, 2013.

Iqbal11 W. Iqbal, M. N. Dailey, D. Carrera, and P. Janecek, "Adaptive Resource Provisioning for Read Intensive Multi-tier Applications in the Cloud," *Future Generation Computer Systems*, vol. 27, no. 6, pp. 871-879, 2011.

Lorido-Botran14 T. Lorido-Botran, J. Miguel-Alonso, and J. A. Lozano, "A Review of Autoscaling Techniques for Elastic Applications in Cloud Environments," *Journal of Grid Computing*, vol. 12, no. 4, pp. 559-592, 2014.

Maurer11 M. Maurer, I. Brandic, and R. Sakellariou, "Enacting Slas in Clouds Using Rules," in *Euro-Par 2011 Parallel Processing*. Springer, 2011, pp. 455-466.

Rao09 J. Rao, X. Bu, C.-Z. Xu, L. Wang, and G. Yin, "VCONF: a Reinforcement Learning Approach to Virtual Machines Auto-configuration," in *Proceedings of the 6th international conference on Autonomic computing*. ACM, 2009, pp. 137-146.

Urgaonkar08 B. Urgaonkar, P. Shenoy, A. Chandra, P. Goyal, and T. Wood, "Agile Dynamic Provisioning of Multi-tier Internet Applications," *ACM Transactions on Autonomous and Adaptive Systems (TAAS)*, vol. 3, no. 1, p. 1, 2008.

Gartner09 D.C. Plume, D. M. Smith, T.J. Bittman, D.W. Cearley, D.J. Cappuccio, D. Scott, R Kumar, and B. Robertson. Study: "Five Refining Attributes of Public and Private Cloud Computing", Tech. rep., Gartner, 2009.

Galante12 G. Galante and L. C. E. d. Bona, "A Survey on Cloud Computing Elasticity" in *Proceedings of the 2012 IEEE/ACM Fifth International Conference on Utility and Cloud Computing*, Washington, 2012

Jennings14 B. Jennings and R. Stadler, "Resource management in clouds: Survey and research challenges", *Journal of Network and Systems Management*, pp. 1-53, 2014

Binning09 C. Binnig, D. Kossmann, T. Kraska, and S. Loesing, "How is the weather tomorrow?: towards a benchmark for the cloud" in *Proceedings of the Second International Workshop on Testing Database Systems*, 2009

Li10 A. Li, X. Yang, S. Kandula, and M. Zhang, "CloudCmp: Comparing Public Cloud Providers" in *Proceedings of the 10th ACM SIGCOMM Conference on Internet Measurement*, 2010

Dory11 T. Dory, B. Mejías, P. V. Roy, and N.-L. Tran, "Measuring Elasticity for Cloud Databases" in *Proceedings of the The Second International Conference on Cloud Computing, GRIDs, and Virtualization*, 2011

Almeida13 R.F. Almeida, F.R.C. Sousa, S. Lifschitz, and J.C. Machado: "On defining metrics for elasticity of cloud databases", *Simpósio Brasileiro de Banco de Dados - SBBD 2013*, <http://www.lbd.dcc.ufmg.br/colecoes/sbbd/2013/0012.pdf>

Weimann11 J. Weinman, "Time is Money: The Value of "On-Demand",", 2011, http://www.joeweinman.com/resources/Joe_Weinman_Time_Is_Money.pdf,

Islam12 S. Islam, K. Lee, A. Fekete, and A. Liu, "How a consumer can measure elasticity for cloud platforms" in *Proceedings of the 3rd ACM/SPEC International Conference on Performance Engineering*, New York, 2012

Folkerts12 E. Folkerts, A. Alexandrov, K. Sachs, A. Iosup, V. Markl, and C. Tosun, "Benchmarking in the Cloud: What It Should, Can, and Cannot Be" in *Selected Topics in Performance Evaluation and Benchmarking*, Berlin Heidelberg, 2012

Moldovan13 D. Moldovan, G. Copil, H.-L. Truong, and S. Dustdar, "MELA: Monitoring and Analyzing Elasticity of Cloud Services," in *IEEE 5th International Conference on Cloud Computing Technology and Science (CloudCom)*, 2013

Tinnefeld14 C. Tinnefeld, D. Taschik, and H. Plattner, "Quantifying the Elasticity of a Database Management System," in *DBKDA 2014, The Sixth International Conference on Advances in Databases, Knowledge, and Data Applications*, 2014

Papadopoulos16 A. V. Papadopoulos, A. Ali-Eldin, K.-E. Årzén, J. Tordsson, and E. Elmroth. 2016. PEAS: A Performance Evaluation Framework for Auto-Scaling Strategies in Cloud Applications. *ACM Trans. Model. Perform. Eval. Comput. Syst.* 1, 4, Article 15 (August 2016), 31 pages. DOI: <http://dx.doi.org/10.1145/2930659>

Massov11 R. v. Massow, R., A. v. Hoorn, W. Hasselbring W. (2011) Performance Simulation of Runtime Reconfigurable Component-Based Software Architectures. In: Crnkovic I., Gruhn V., Book M. (eds) *Software Architecture*. ECSA 2011. Lecture Notes in Computer Science, vol 6903. Springer, Berlin, Heidelberg

Huber17 N. Huber, F. Brosig, S. Spinner, S. Kounev, and M. Bähr. Model-Based Self-Aware Performance and Resource Management Using the Descartes Modeling Language. *IEEE Transactions on Software Engineering (TSE)*, 43(5), 2017, IEEE Computer Society.