

# Diagnosis and Automation for Deployed Distributed Systems

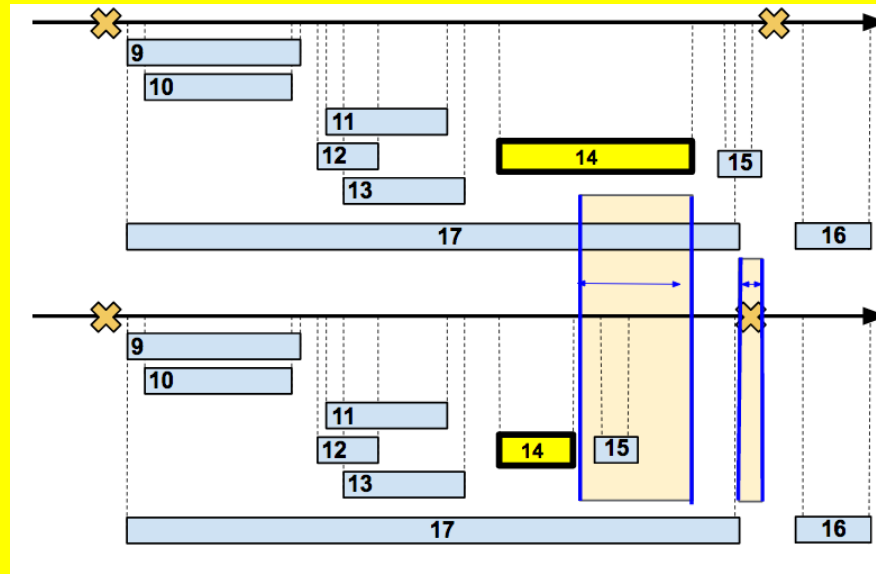
Google™



## The Weatherman Effort

Gideon Mann, Mark Sandler, Eyal Even-dar, Darja Krushevskaja, Sudipto Guha, Krzysztof Ostrowski, Sebastian Pueblas, Lev Ratinov, Eran Gabber

# Today: Profiling Latency in Deployed Distributed Systems



ation  
ed



## The Weatherman Effort

Gideon Mann, Mark Sandler, Eyal Even-dar, Darja Krushevskaja, Sudipto Guha, Krzysztof Ostrowski, Sebastian Pueblas, Lev Ratinov, Eran Gabber

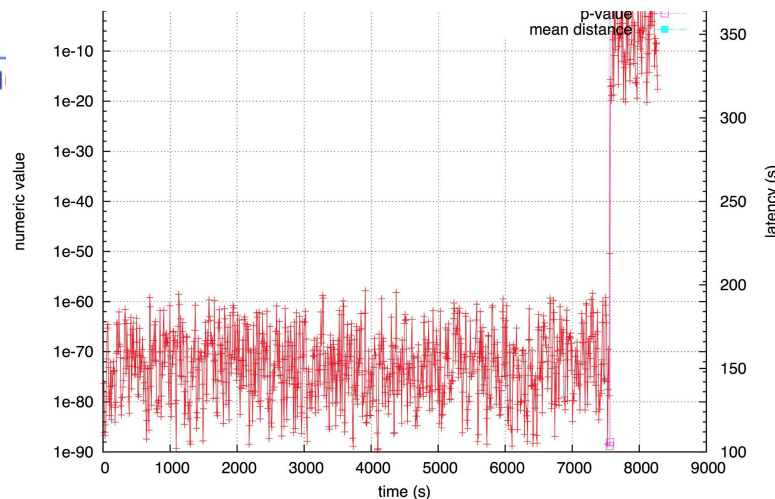
# Two Scenarios of Interest

Latency distribution has multiple peaks, at different orders of magnitude.



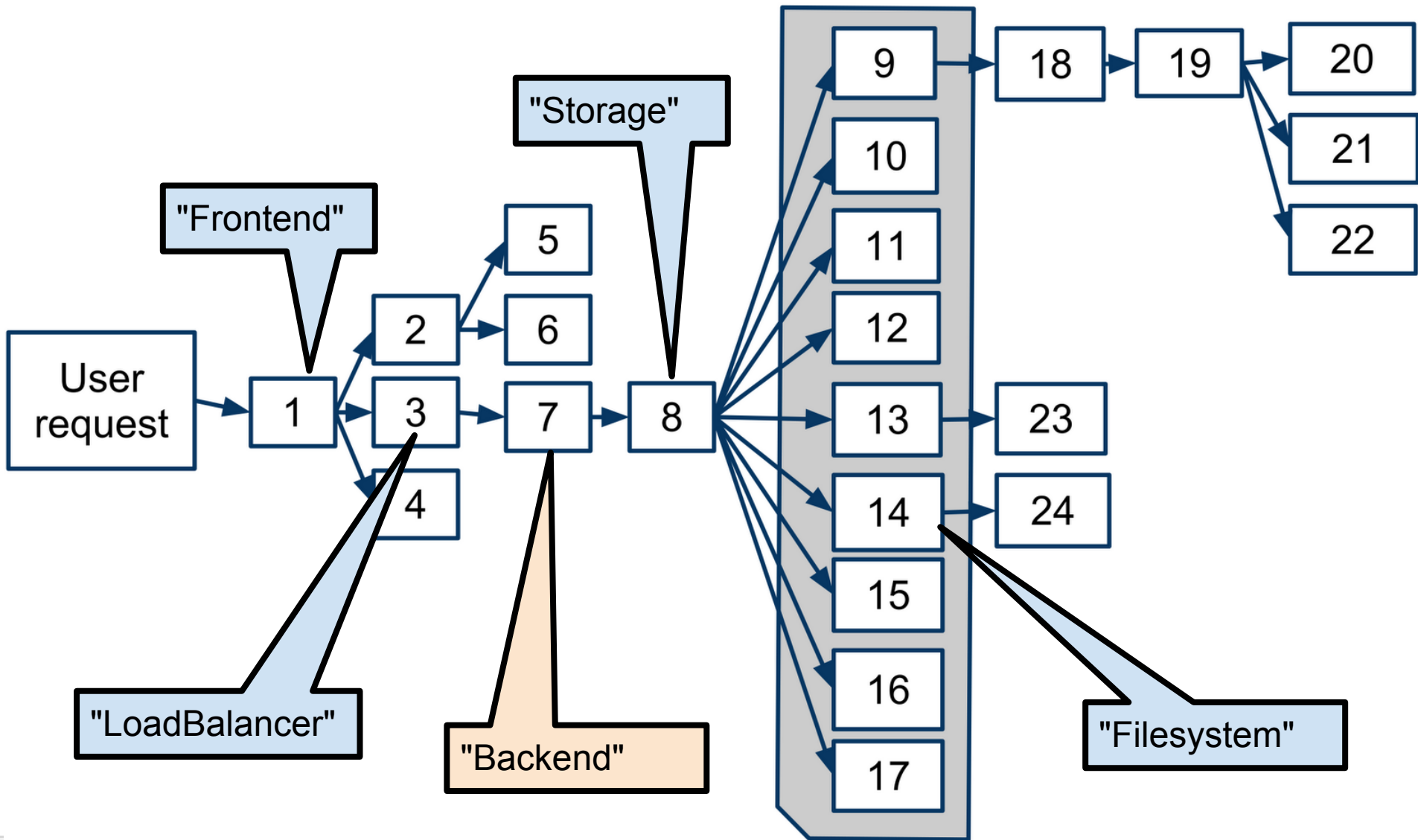
Over time, the latency has dramatically changed ... what happened?

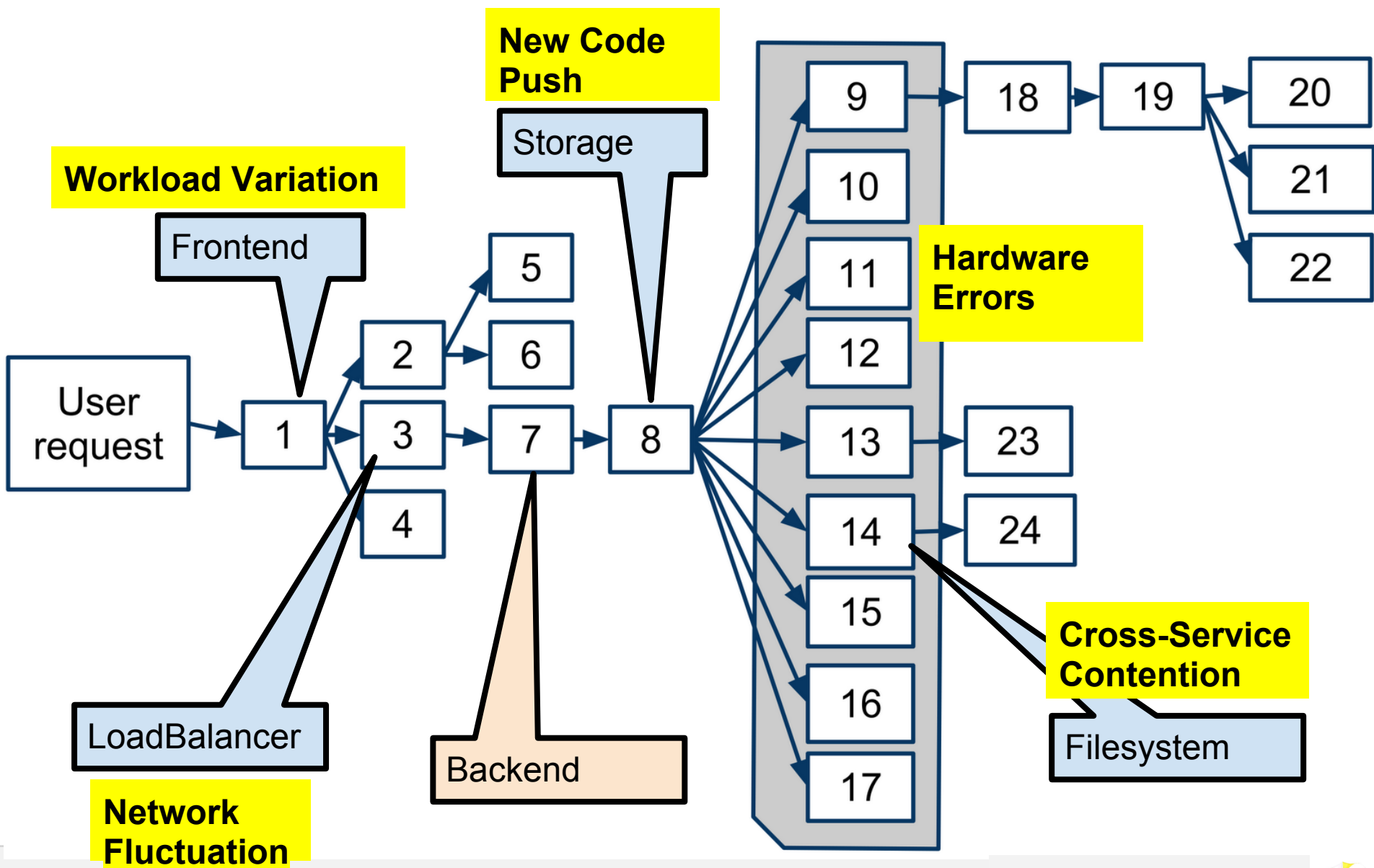
Load-testing isn't enough.



Well, probably a lot has changed. What's the most significant?

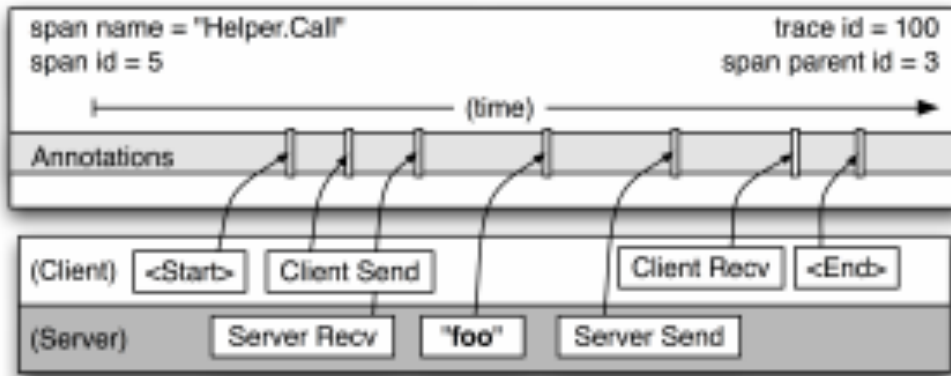
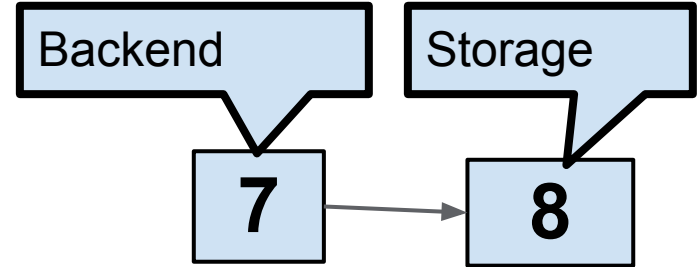






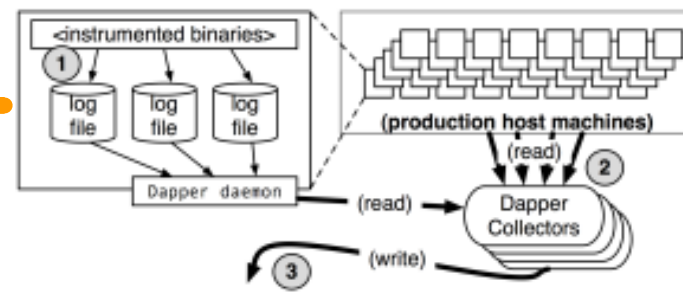
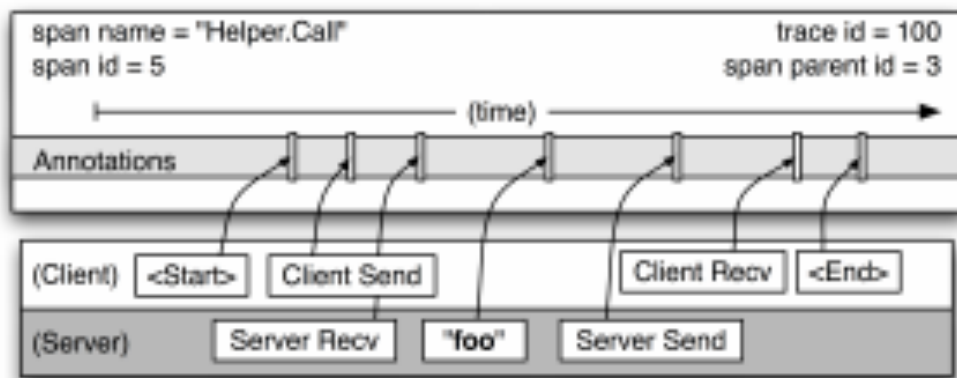
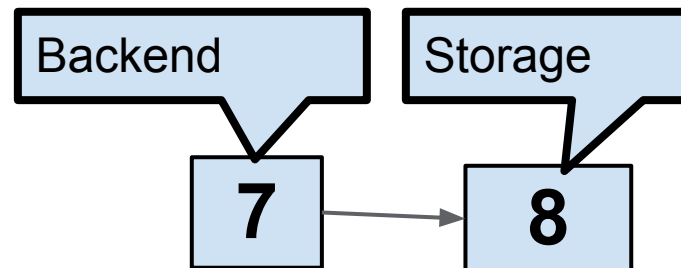
# Dapper Distributed Trace Collection

Sigelman, Barroso, Burrows, Stephenson,  
Plakal, Beaver, Jaspan, Shanbhag '10



# Dapper Distributed Trace Collection

Sigelman, Barroso, Burrows, Stephenson, Plakal, Beaver, Jaspan, Shanbhag '10

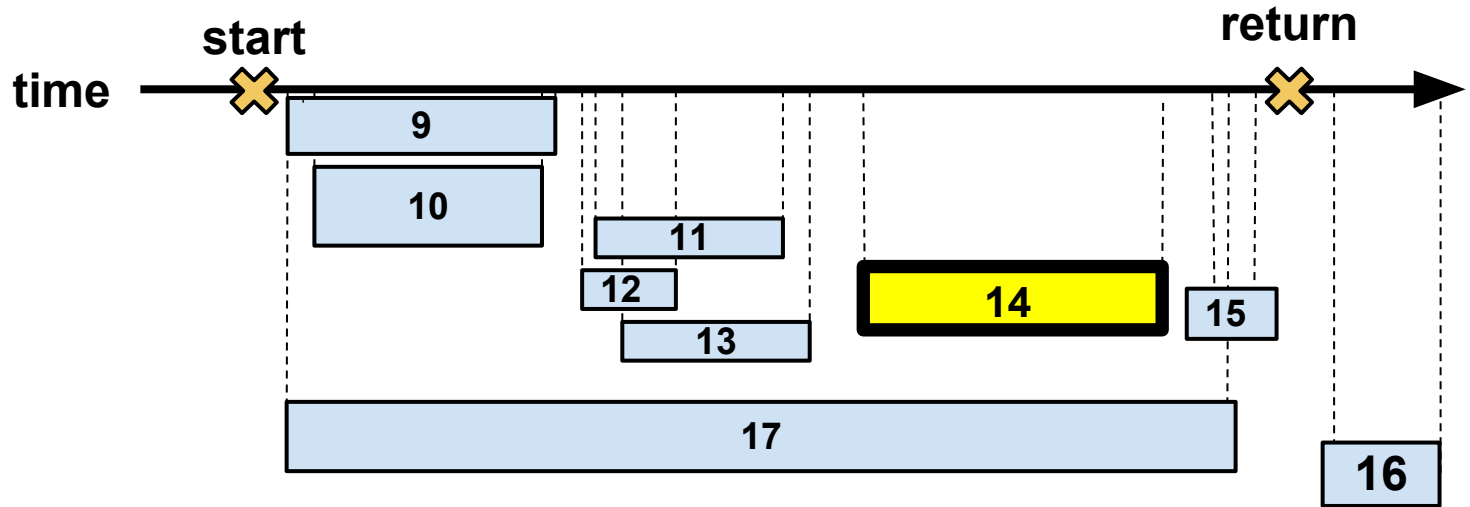
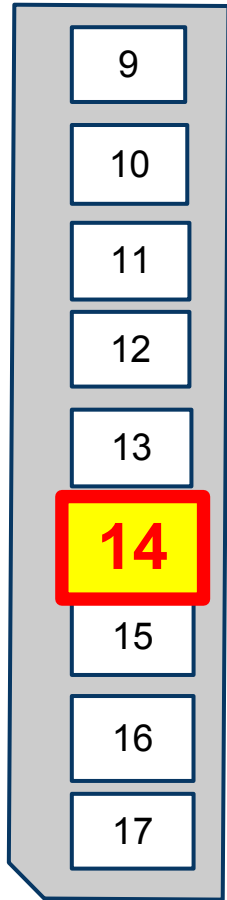
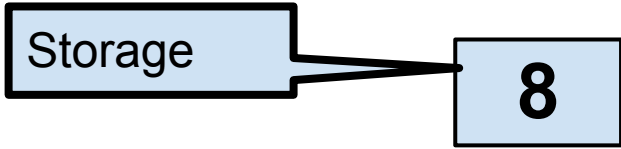


trace id	span 12	span 23	span 34	span 45	span 56	...
123456	nil	nil	<data>	<data>	nil	...
246802	<data>	nil	nil	nil	<data>	...
357913	nil	<data>	nil	nil	nil	...
...	...	...	...	...	...	...

(Central Bigtable repository for trace data)

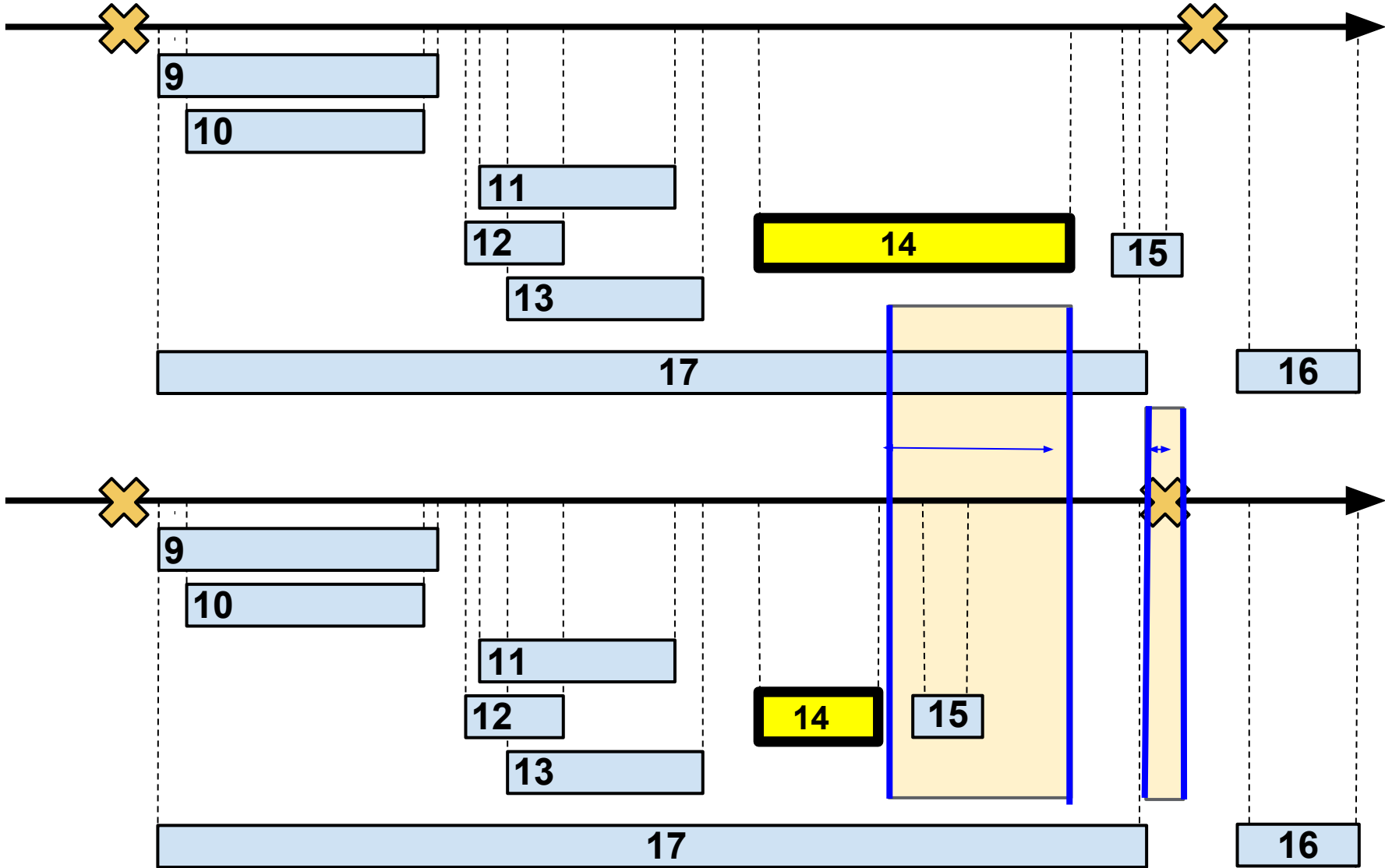
Figure 5: An overview of the Dapper collection pipeline.

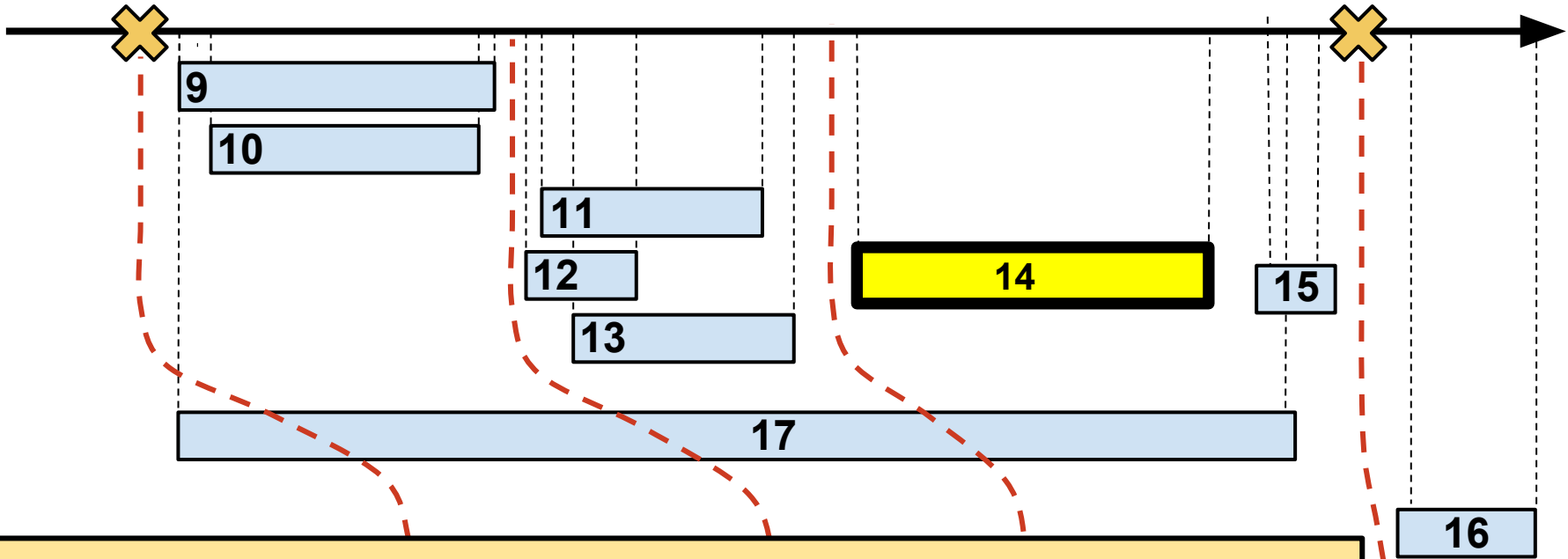




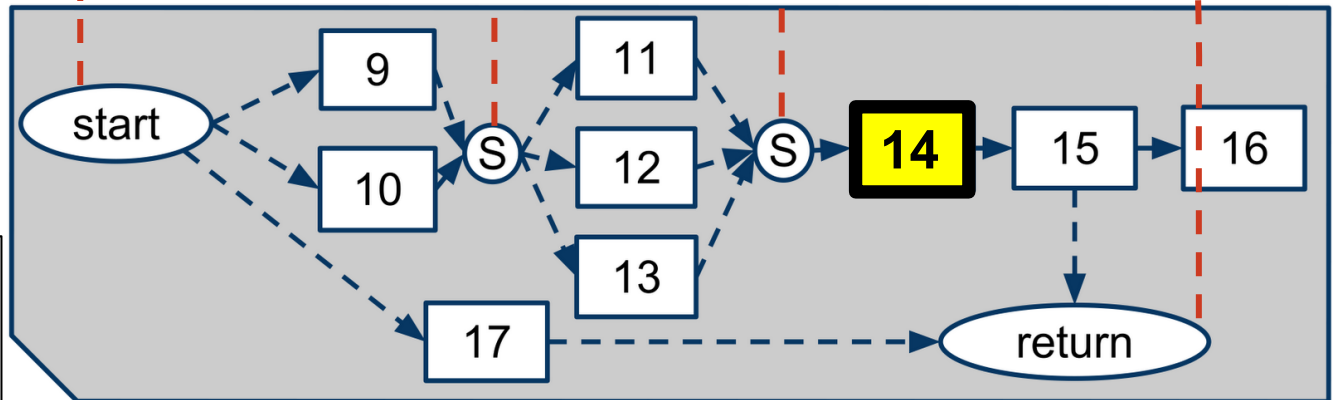


"If time in 14 is reduced, what is the predicted latency?"





We need to model the asynchrony to understand the latency.



Mann, Sandler, Kruschevskaja, Guha, Even-dar. HotCloud '10.

$\lambda_t$       observed latency for RPC  $t$

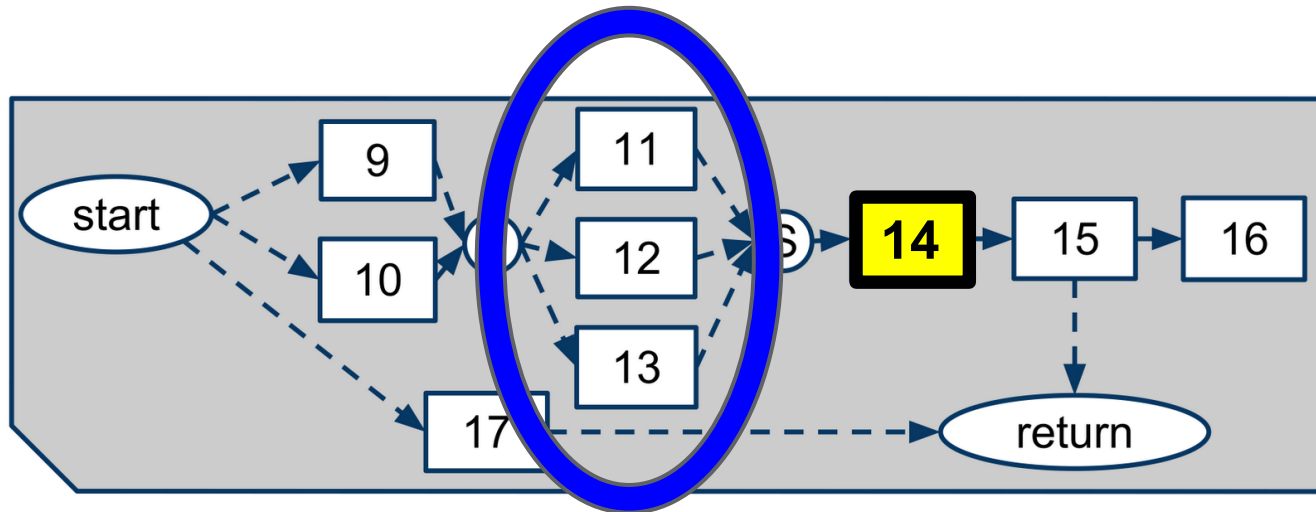
$\hat{\lambda}_t$       predicted latency for RPC  $t$

We predict the latency of a given RPC given the latency of its children  $C(t)$ , and relevant features of the RPC  $t$  and children.

$$\hat{\lambda}_t = \Psi(f_t, f_{C(t)}, \lambda_{C(t)})$$

$\chi$  An (asynchronous) execution model

$\chi(t_i)$  The blocking set for a RPC  $t_i$



$\chi(t_{14})$

$\chi$  An (asynchronous) execution model

$\chi(t_i)$  The blocking set for a RPC  $t_i$

$$\psi_{\chi(t_i)} = \lambda_{t_i} + \max_{t \in \chi(t_i)} \psi(t)$$

The finishing time for a RPC  $t_i$  is its own latency and the finishing times of the RPCs in its blocking set.

$$\begin{aligned} \hat{\lambda}_p &= \Psi(f_p, f_{C(p)}, \lambda_{C(p)}) \\ &= \Psi(\chi, \lambda_{C(p)}) \\ &= \psi_{\chi}(t_{return}) \end{aligned}$$

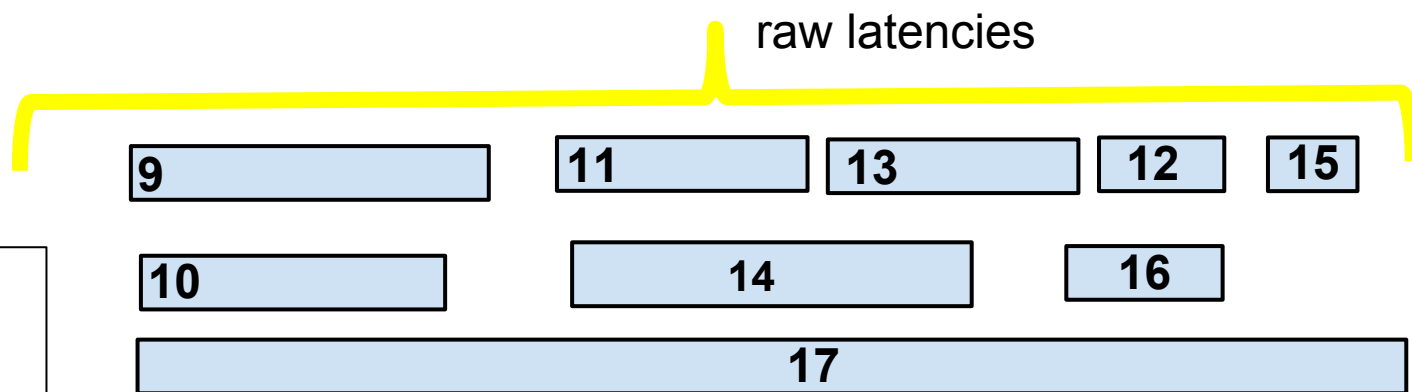
Given the children latencies and the execution model, we can compute the latency for a given parent RPC.



# Find best matching flow via Nearest Neighbor :

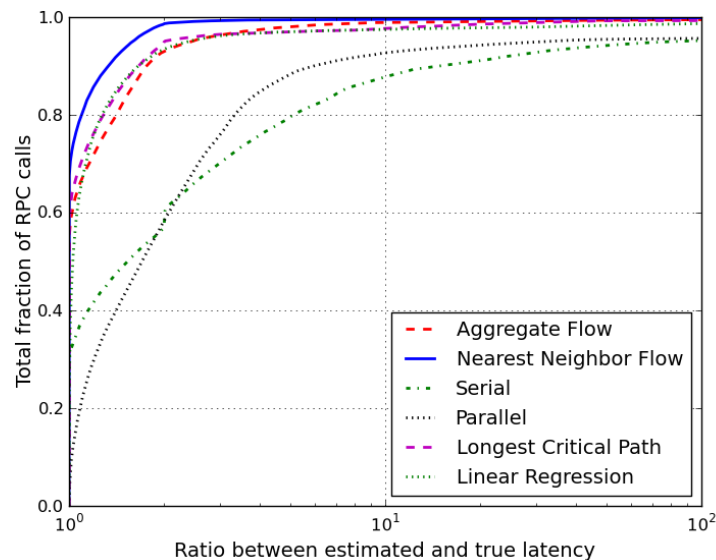
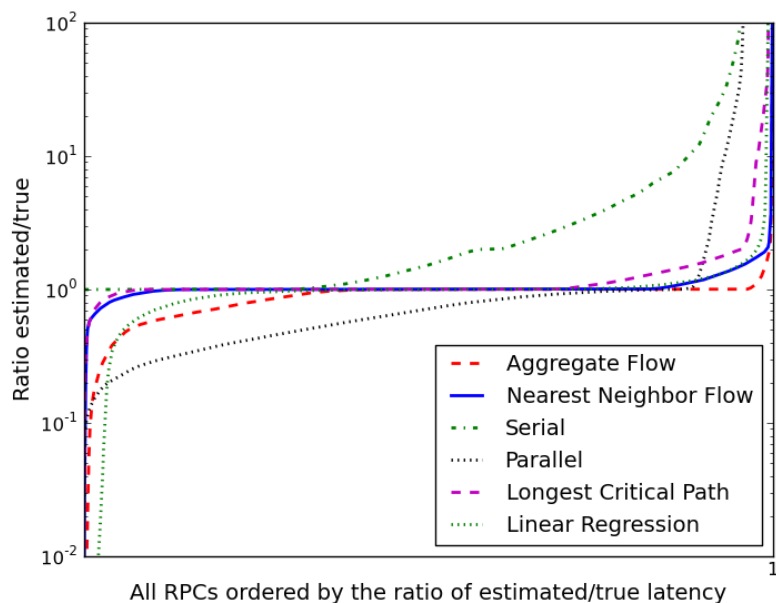
$$\lambda_{t_i}^x \sim N(\bar{\lambda}_{t_i}, \bar{\sigma}_{t_i})$$

$$\log P_x(\lambda_t) \propto \sum_{t_i \in C(t)} (\bar{\lambda}_{t_i} - \lambda_{t_i}) / \bar{\sigma}_{t_i}$$



Mann, Sandler,  
Kruschevskaja,  
Guha, Even-dar.  
HotCloud '10.

A perfect model would be a line at  $y=1.0$



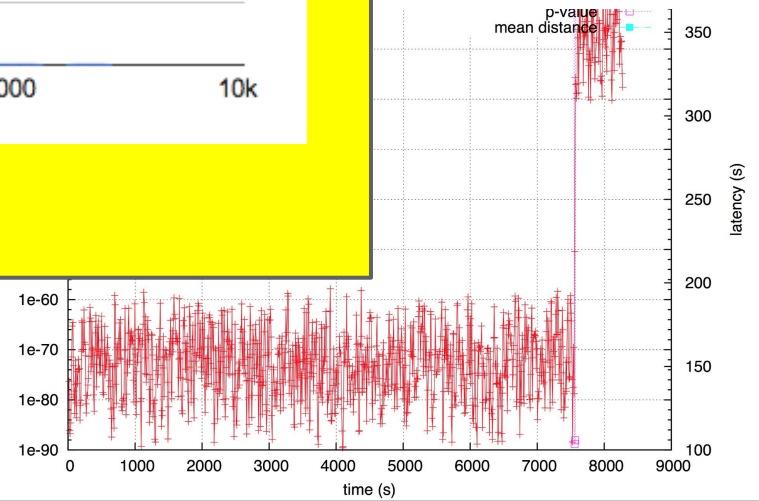
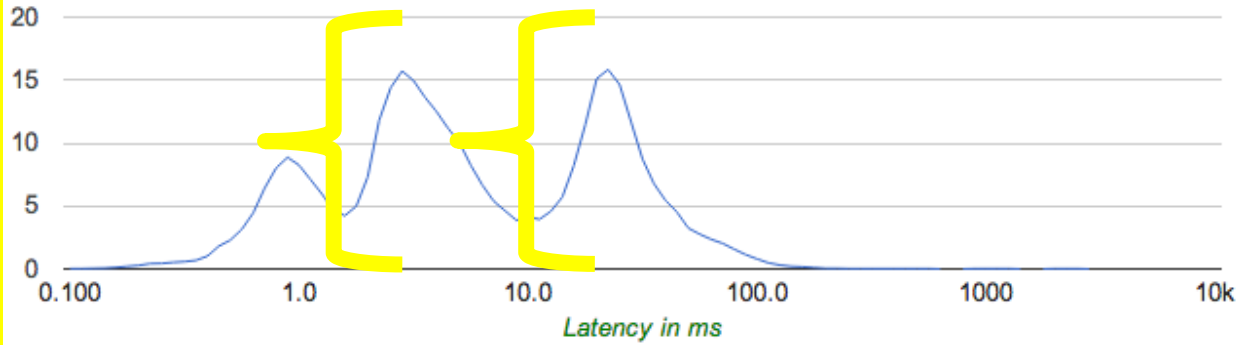
Induced models ("Nearest Neighbor Flow") are better predictors of parent latencies from children than:

- Linear regression
- Longest critical path
- Either a full serial or full parallel RPCs

Mann, Sandler,  
Kruschevskaja,  
Guha, Even-dar.  
HotCloud '10.

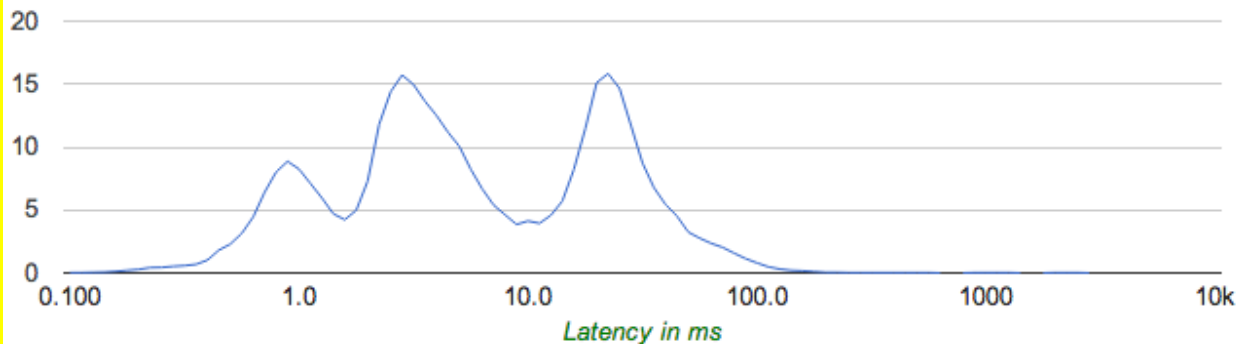
# First Scenario of Interest

Latency distribution has multiple peaks at different orders of magnitude.





Latency distribution has multiple peaks at different orders of magnitude.

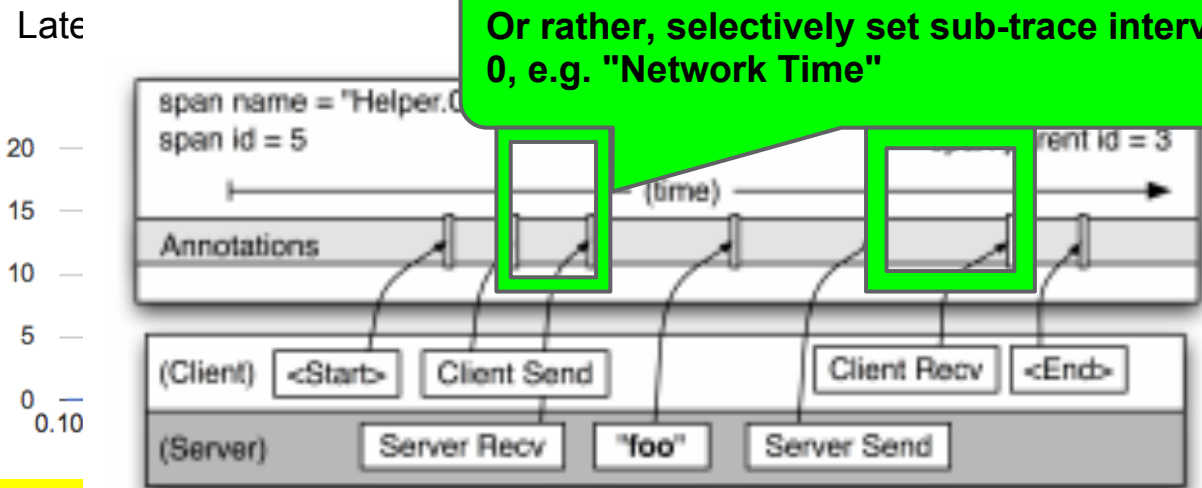


$$\Lambda_A = \bigcup_{t \in A} (\lambda_t)$$

$$\Lambda_{\delta A} = \mathbb{P}_{t \in A}(\Psi(\chi, \lambda_{\delta C(t)}))$$

Selectively set latencies of specific traces to 0, and then simulate the effects.

Or rather, selectively set sub-trace intervals to 0, e.g. "Network Time"



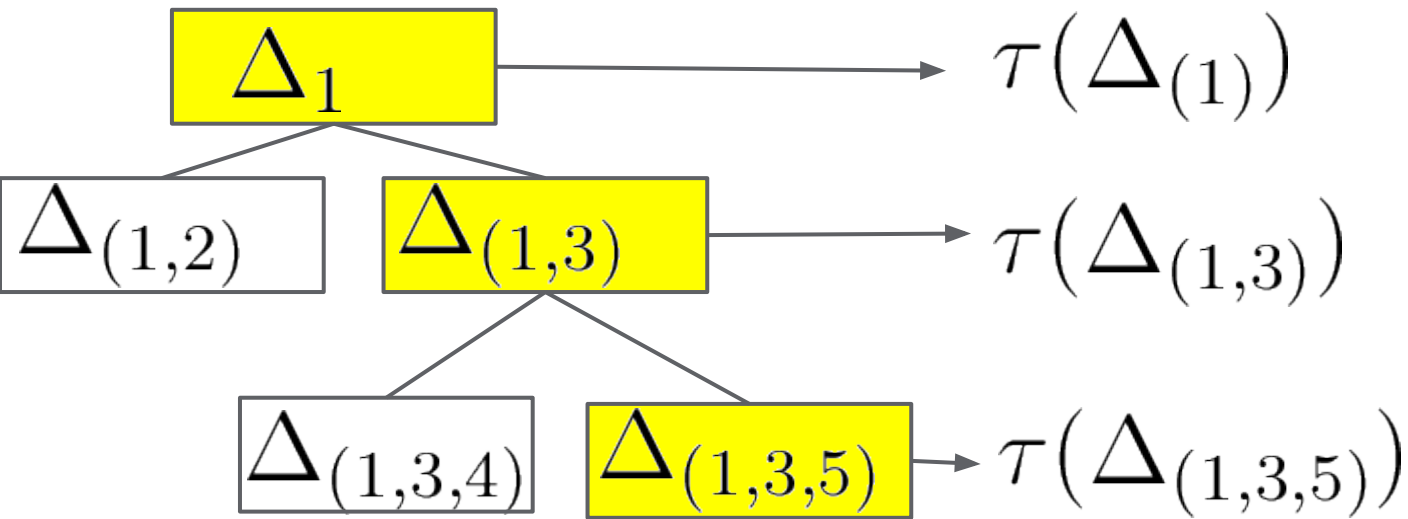
$$\Lambda_A = \bigcup_{t \in A} (\lambda_t)$$

$$\Lambda_{\delta A} = \mathbb{P}_{t \in A} (\Psi(\chi, \lambda_{\delta C(t)}))$$

Selectively set latencies of specific traces to 0, and then simulate the effects.

$\Lambda_A$  Random variable of Latency

$$\tau(\Delta_j) = \mathbb{E}[\Lambda_{\Delta_j A}]$$



$\Lambda_A$  Random variable of Latency

$$\tau(\Delta_j) = \mathbb{E}[\Lambda_{\Delta_j A}]$$

Mean Latency at 90% percentile

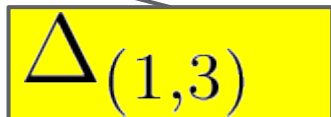
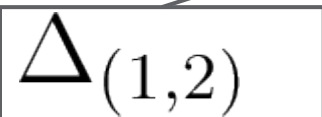
100ms



Total Bigtable Latency -> 0

$$\tau(\Delta_{(1)})$$

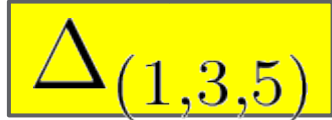
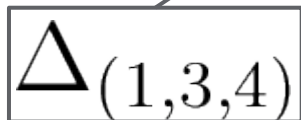
30ms



Bigtable Local Latency -> 0

$$\tau(\Delta_{(1,3)})$$

44ms



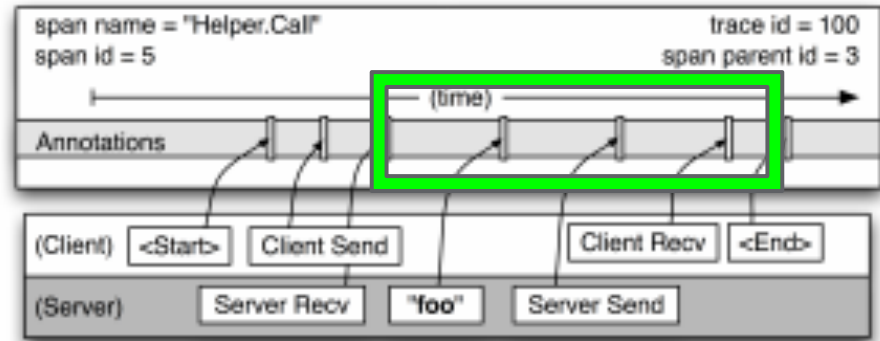
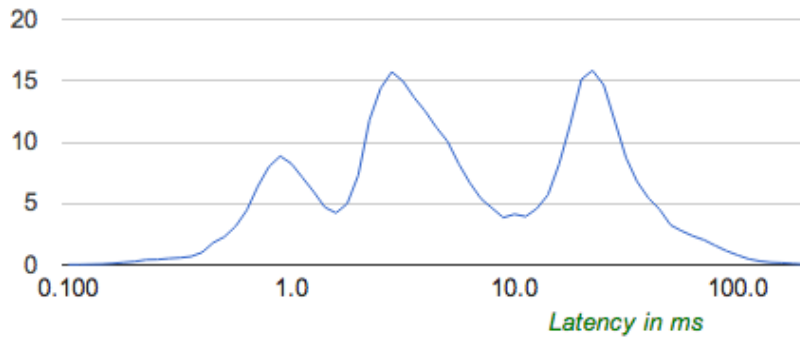
If response payload >10k -> 0

$$\tau(\Delta_{(1,3,5)})$$

45ms

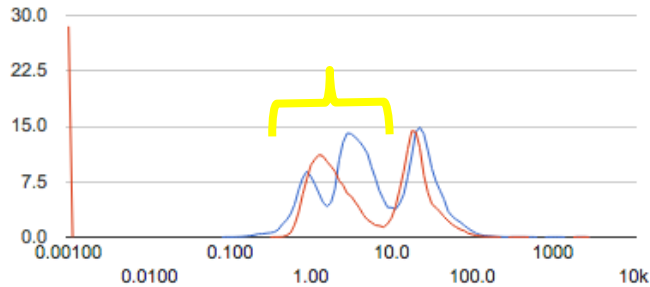


### Overall latency distribution

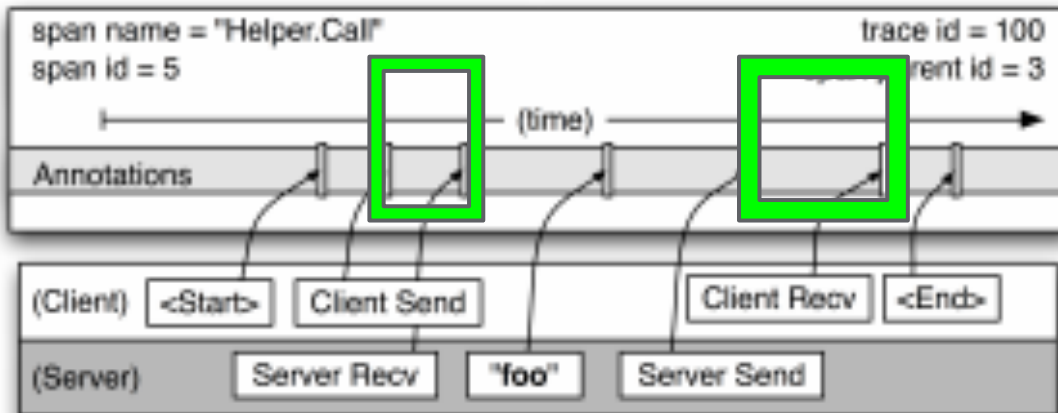
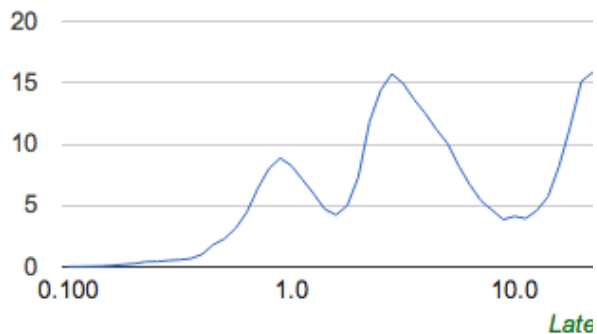


Red Line : top-level latency without time in ...

### Bigtable Local Processing

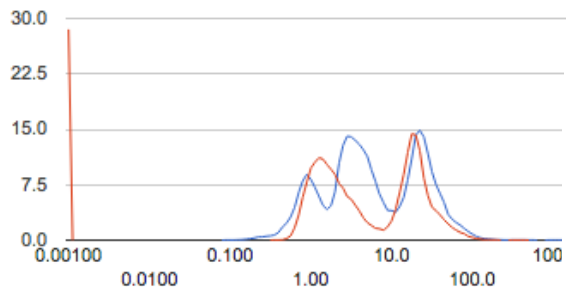


### Overall latency distribution

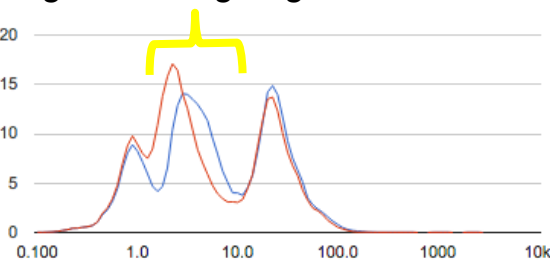


Red Line : top-level latency with

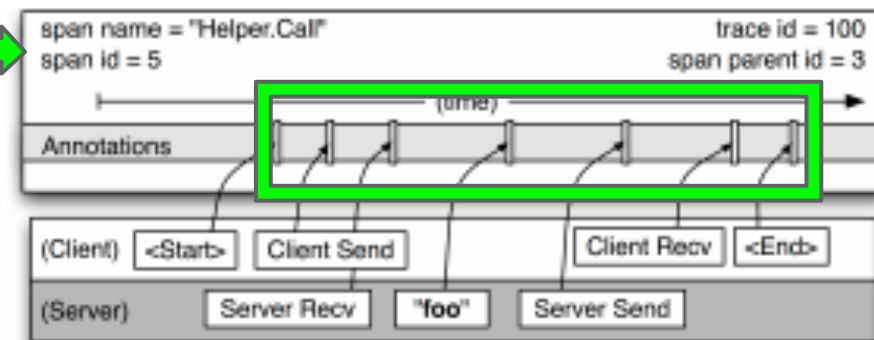
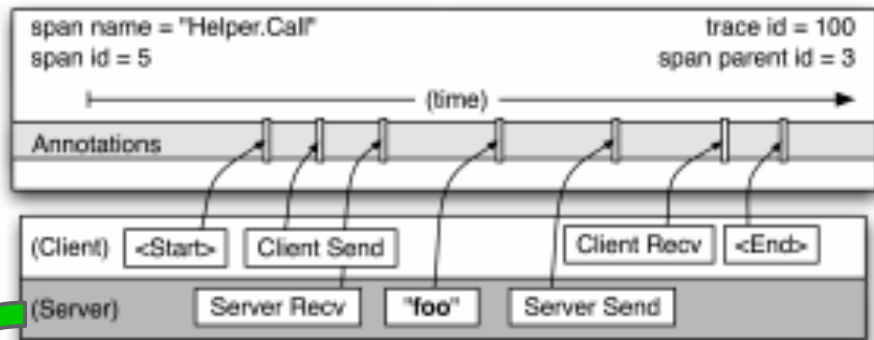
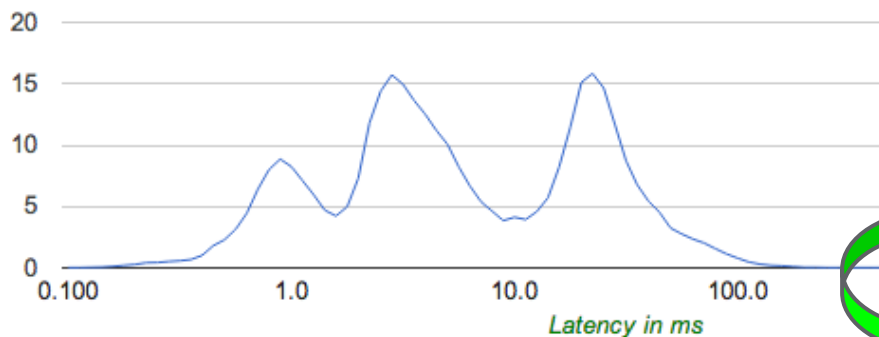
### Bigtable Local Processing



### Bigtable Outgoing Network

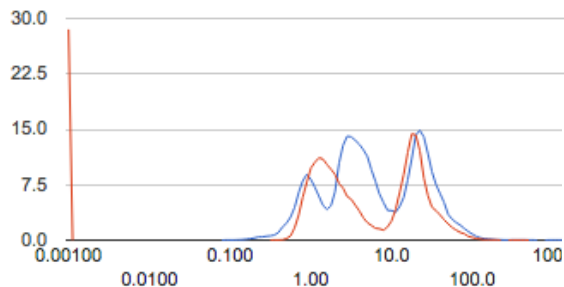


Overall latency distribution

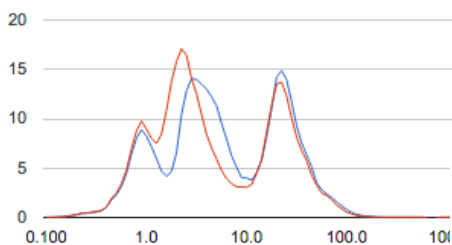


Red Line : top-level latency without time in ...

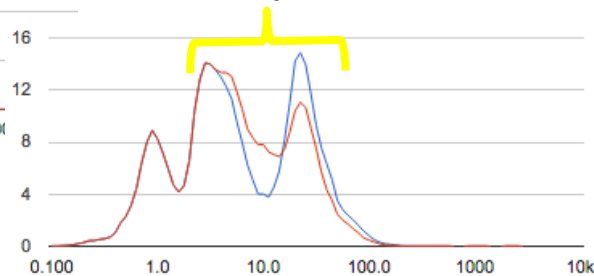
Bigtable Local Processing



Bigtable Outgain

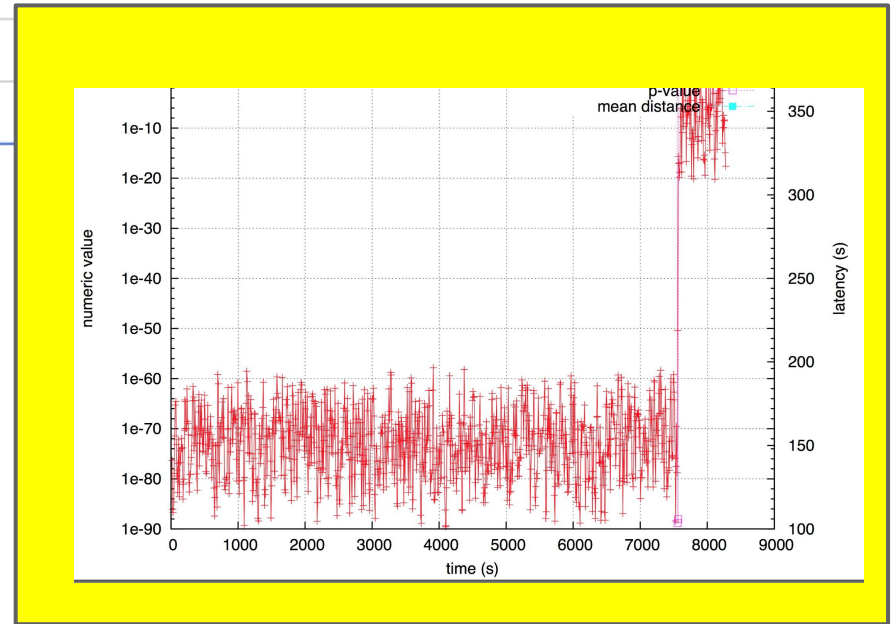
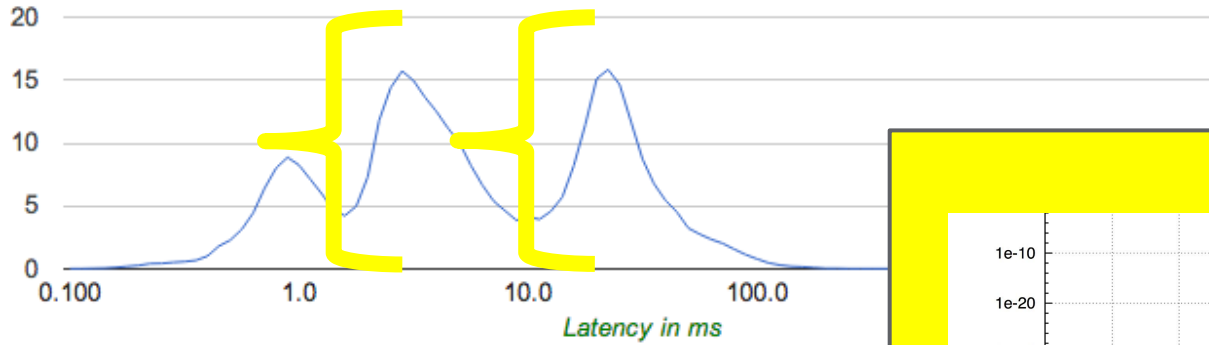


Total FileSystem Time

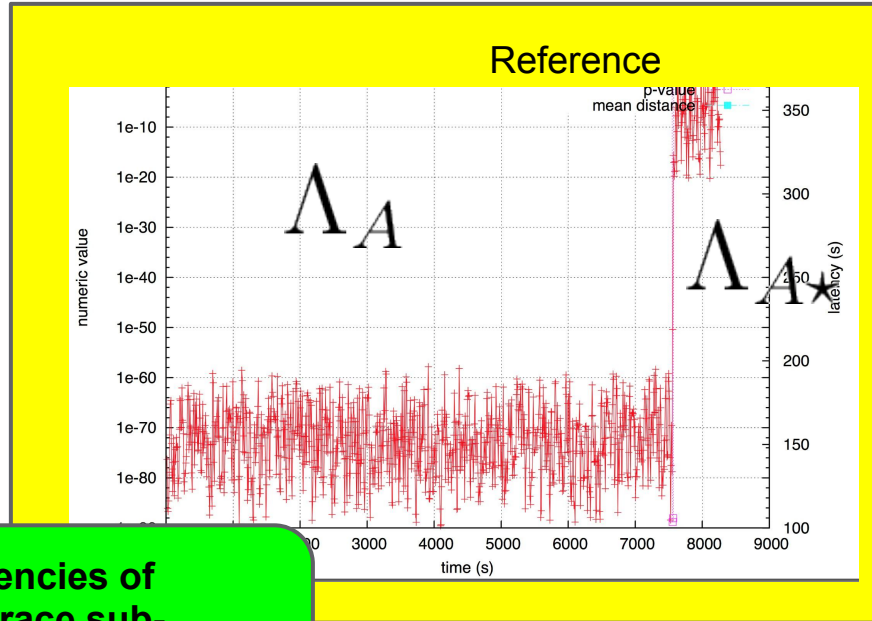


# Second Scenario of Interest

Latency distribution has multiple peaks at different orders of magnitude.







Selectively alter latencies of specific traces (or trace sub-components) and then simulate the effects.

$$\Lambda_{\delta A} = \mathbb{P}_{t \in A}(\Psi(\chi, \lambda_{\delta C}(t)))$$

$$\Lambda_{\delta A} = \mathbb{P}_{t \in A}^{\delta}(\lambda_t)$$

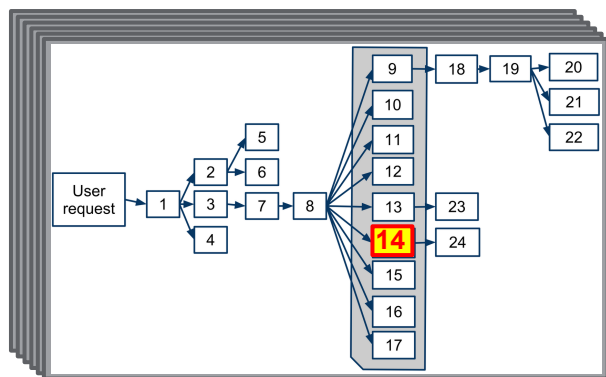
Alter overall distribution of traces and then simulate the effects.

# Sample Based Approach

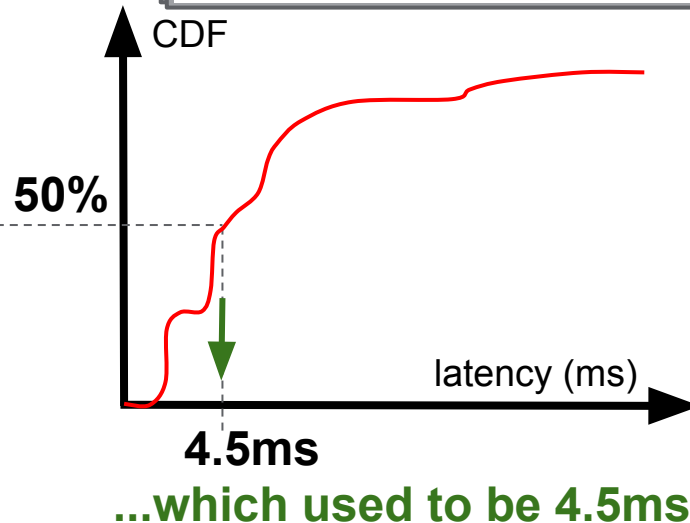
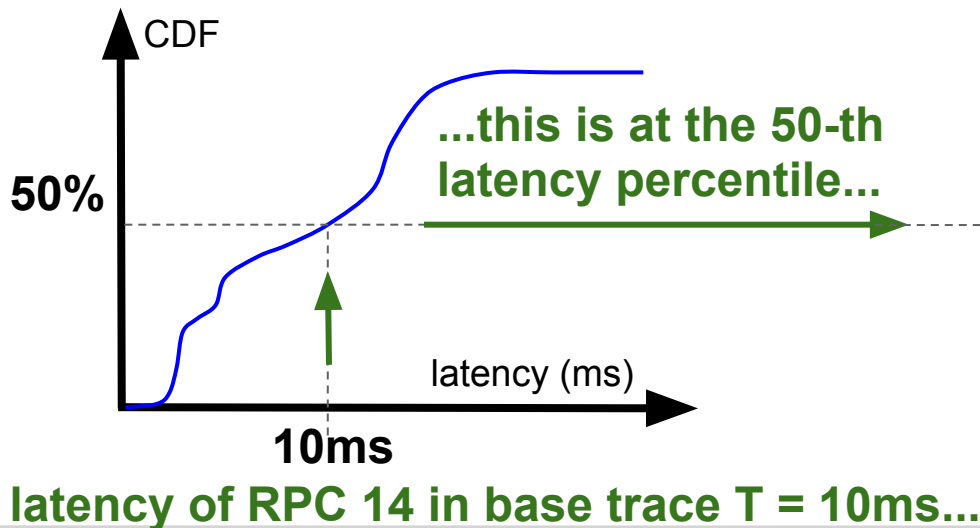
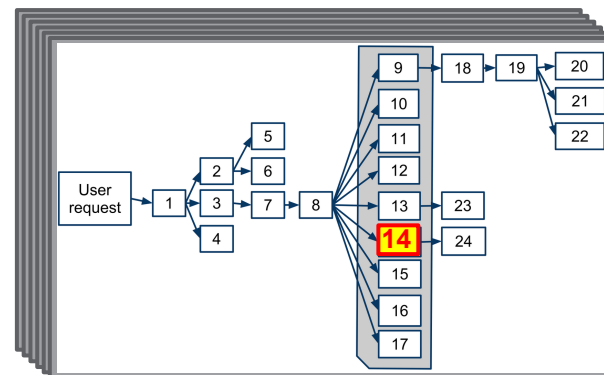
=> shift latency distribution

$$\Lambda_{\delta A} = \mathbb{P}_{t \in A}(\Psi(\chi, \lambda_{\delta C(t)}))$$

"base" traces



"target" traces



# Sample Based Approach

=> shift distribution of traces

$$\Lambda_{\delta A} = \mathbb{P}_{t \in A}^{\delta}(\lambda_t)$$

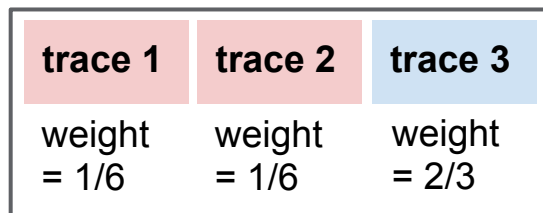
= absent  
 = present

"base" traces



"fix"

modified traces



"target" traces



Root Cause	25%	50%	75%	90%	95%	98%
▶ Total Server <a href="#">Detail</a>			-8.67	-1048.20	-1618.00	-583.34
▶ ./Storage Read <a href="#">Detail</a>	0.00	0.02	-9.37	-1015.80	-1516.00	-376.00
▶ Total Local Computation <a href="#">Detail</a>	0.00	0.02	-9.36	-1015.80	-1516.00	-376.00
▶ :--method returned-- <a href="#">Detail</a>	0.00	0.02	-9.36	-1015.80	-1516.00	-376.00
▶ Distribution of Annotation Existence: netsched send (624) <a href="#">Detail</a> @Annotation Existence: netsched send (624): PRESENT <a href="#">Detail</a> Distribution of Annotation Existence: netsched queuing send behind %d others (624) <a href="#">Detail</a> @Annotation Existence: netsched queuing send behind %d others (624): PRESENT <a href="#">Detail</a>	0.00	0.01	-1.89	-950.00	-1501.00	-382.29
▶ @Annotation Existence: starting lock acquisition (610): PRESENT <a href="#">Detail</a>	-0.33	-0.47	-9.85	-1034.20	-1850.00	-610.00

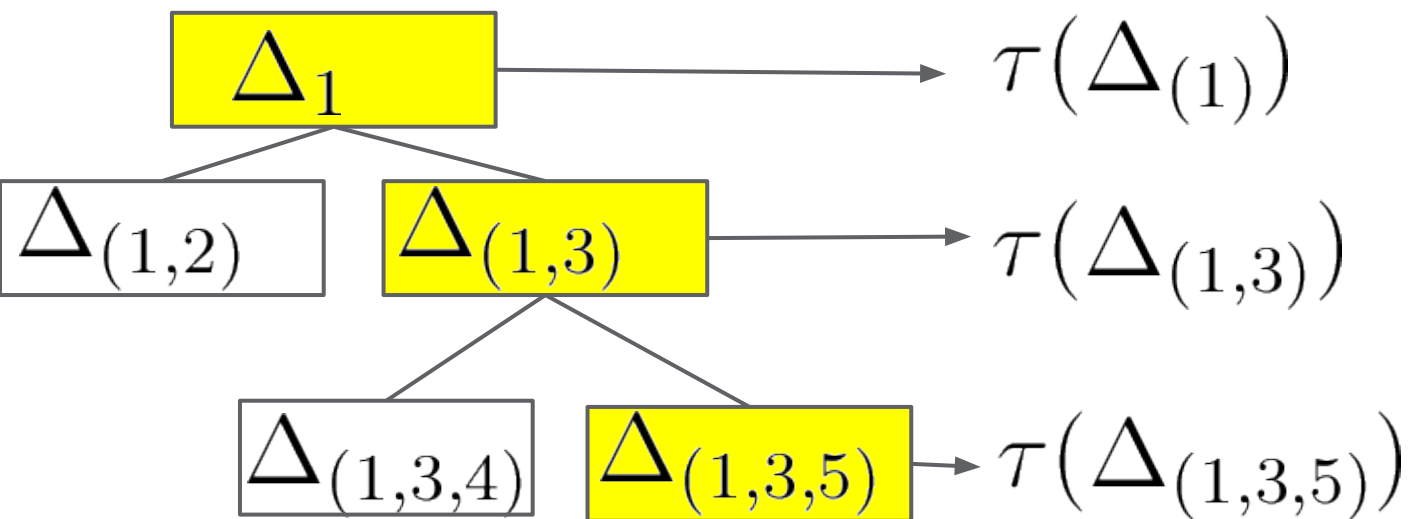
Sample Expectation at of 90%.

Distribution of Annotation Existence: netsched send (624) [Detail](#)  
 @Annotation Existence: netsched send (624): PRESENT [Detail](#)  
 Distribution of Annotation Existence: netsched queuing send behind %d others (624) [Detail](#)  
 @Annotation Existence: netsched queuing send behind %d others (624): PRESENT [Detail](#)



$\Lambda_A$  Sample       $\Lambda_{A^*}$  Reference

$$\tau(\Delta_j) = \mathbb{E}[\Lambda_{\Delta_j A}] - \mathbb{E}[\Lambda_{A^*}]$$

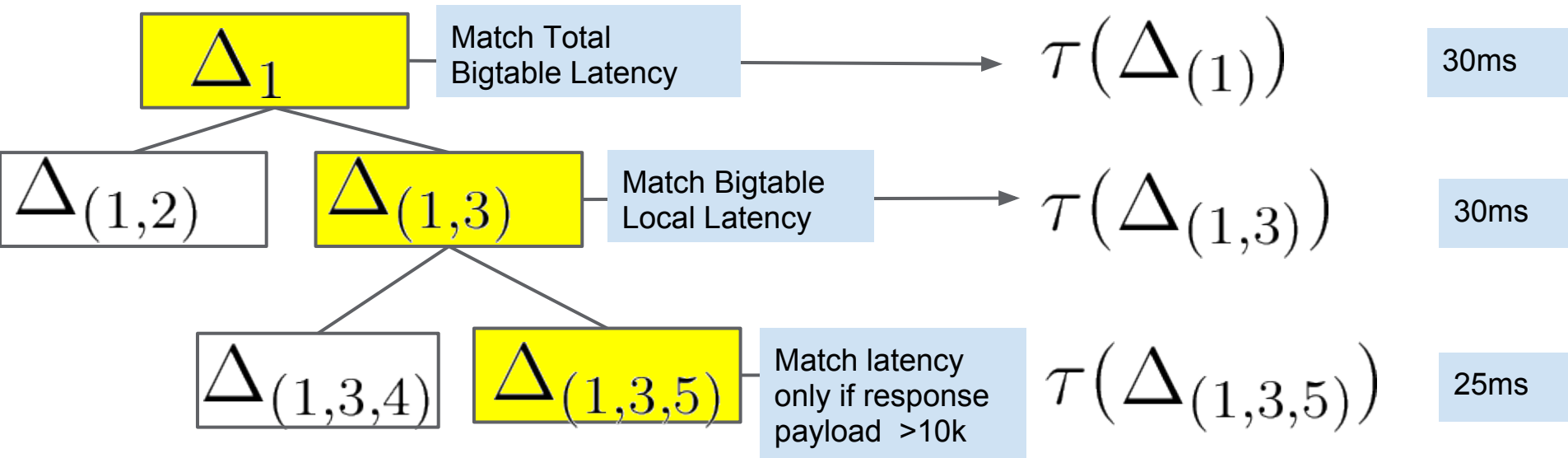


$\Lambda_A$  Sample       $\Lambda_{A^*}$  Reference

$$\tau(\Delta_j) = \mathbb{E}[\Lambda_{\Delta_j A}] - \mathbb{E}[\Lambda_{A^*}]$$

Mean Latency  
Difference at  
90% percentile

50ms

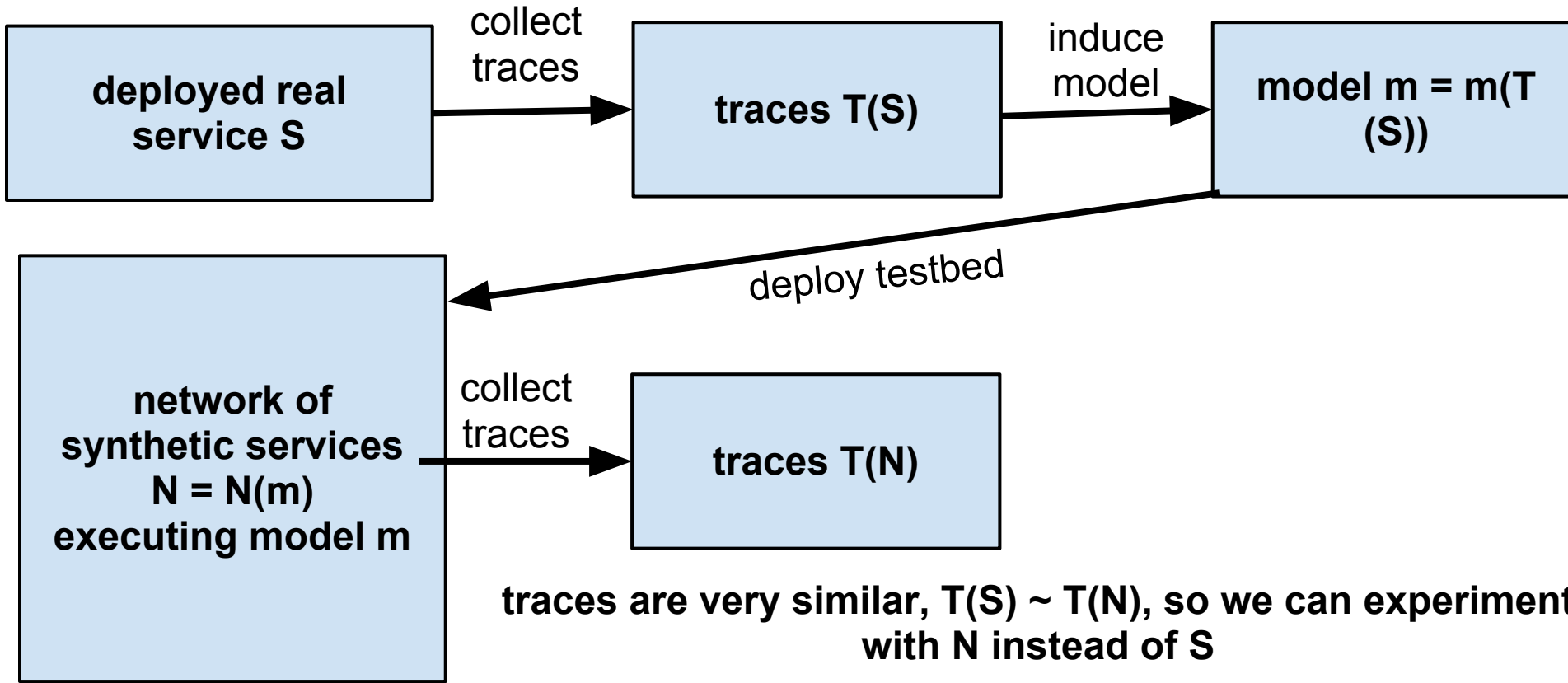


# Evaluation Testbed

**Synthetic Services:** physical network delays, random lognormal processing times, temporal dependencies from the induced execution models

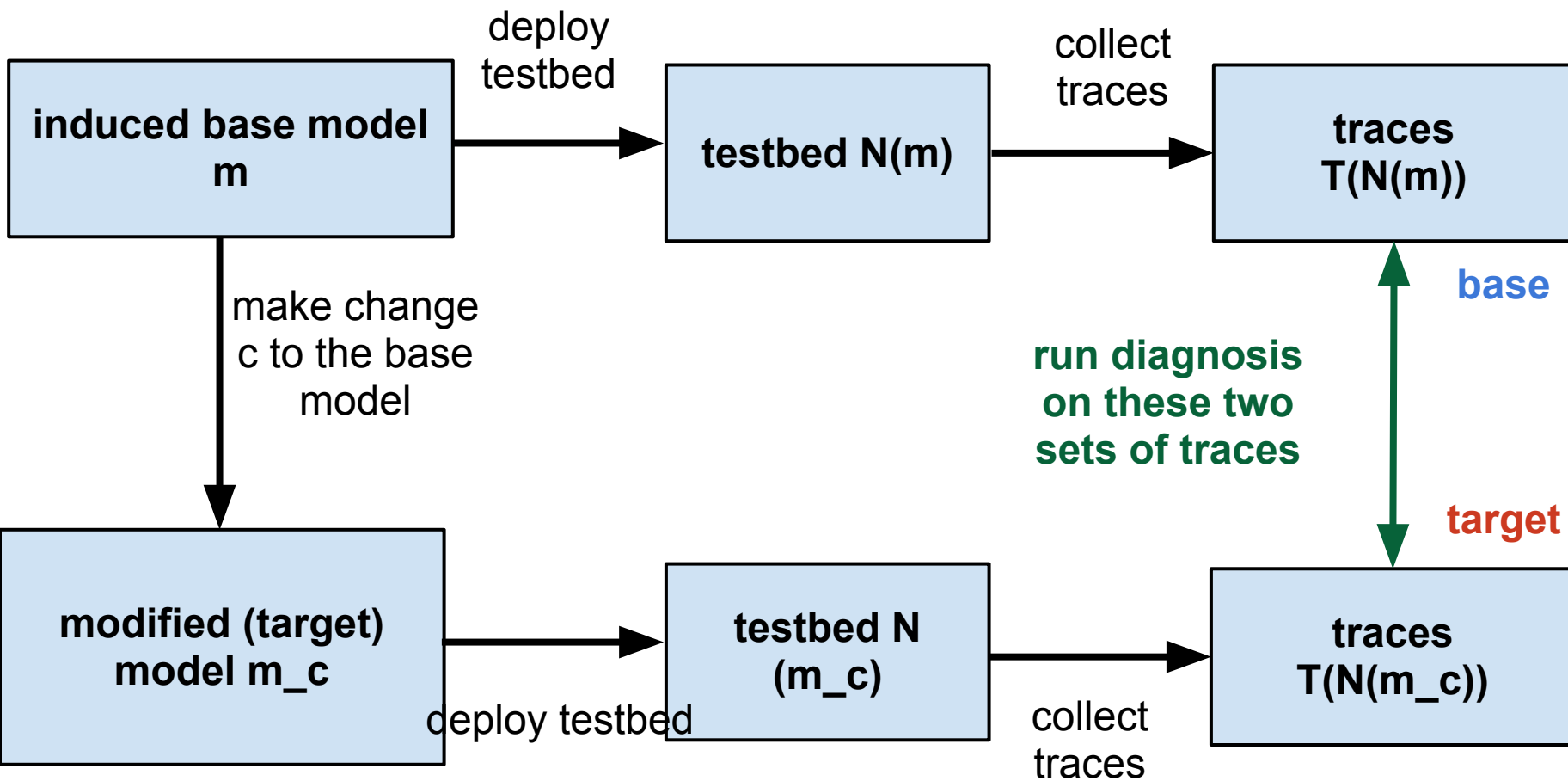
**Diagnose changes we introduce:** Modify the model of a service, generate a new set of traces using the modified model, and diagnose the difference between these sets of traces using our tool





Ostrowski, Mann, Sandler. LADIS '11.





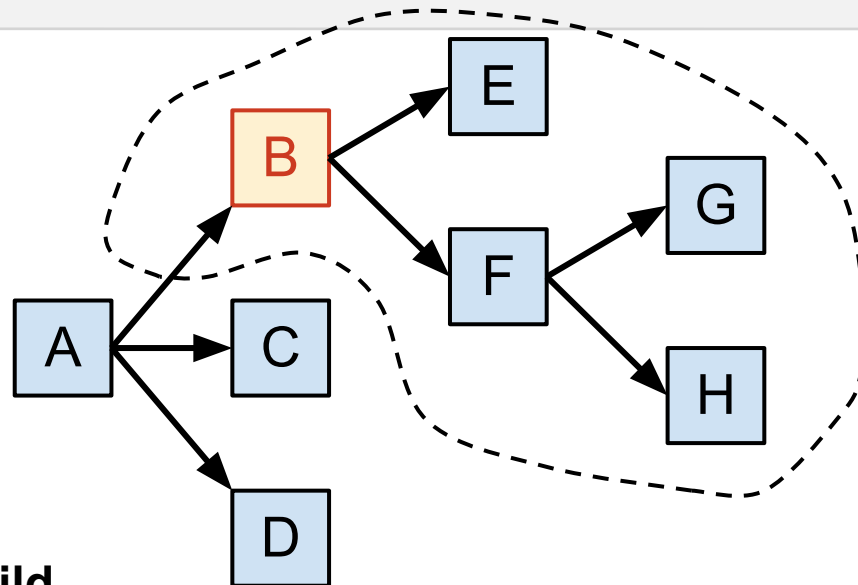
Ostrowski, Mann, Sandler. LADIS '11.



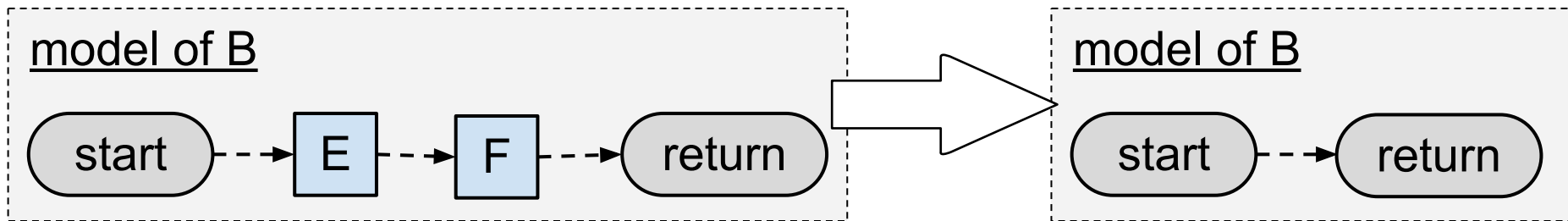


## Two Introduced Changes

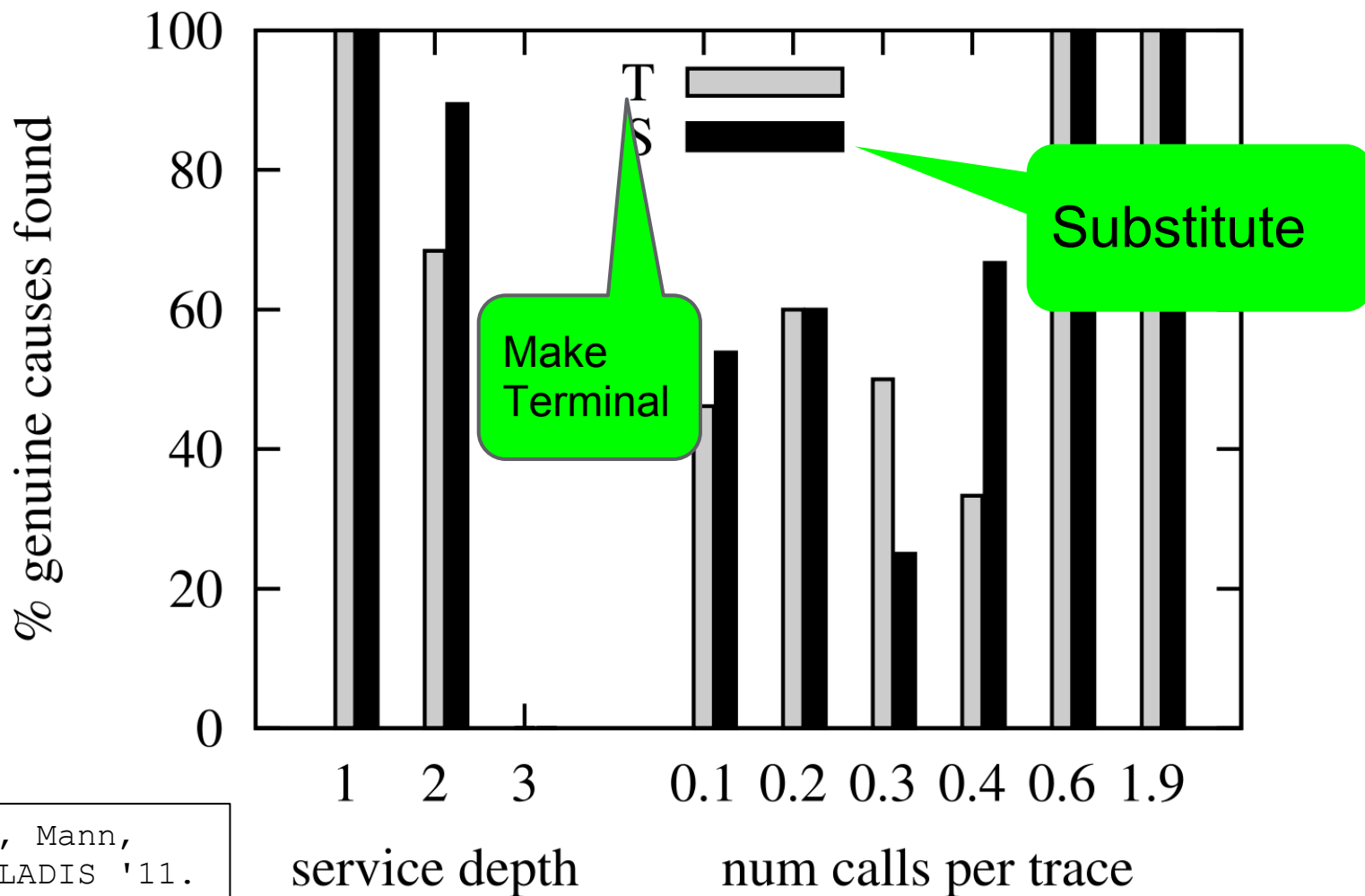
**1. Substitute ("S"):** replace the models in a subtree of base services with those from the target. Only on one call path, so if B is called in two different places, only one of them is affected



**2. Make Terminal ("T"):** remove all child RPCs in a model of a service; make it return immediately

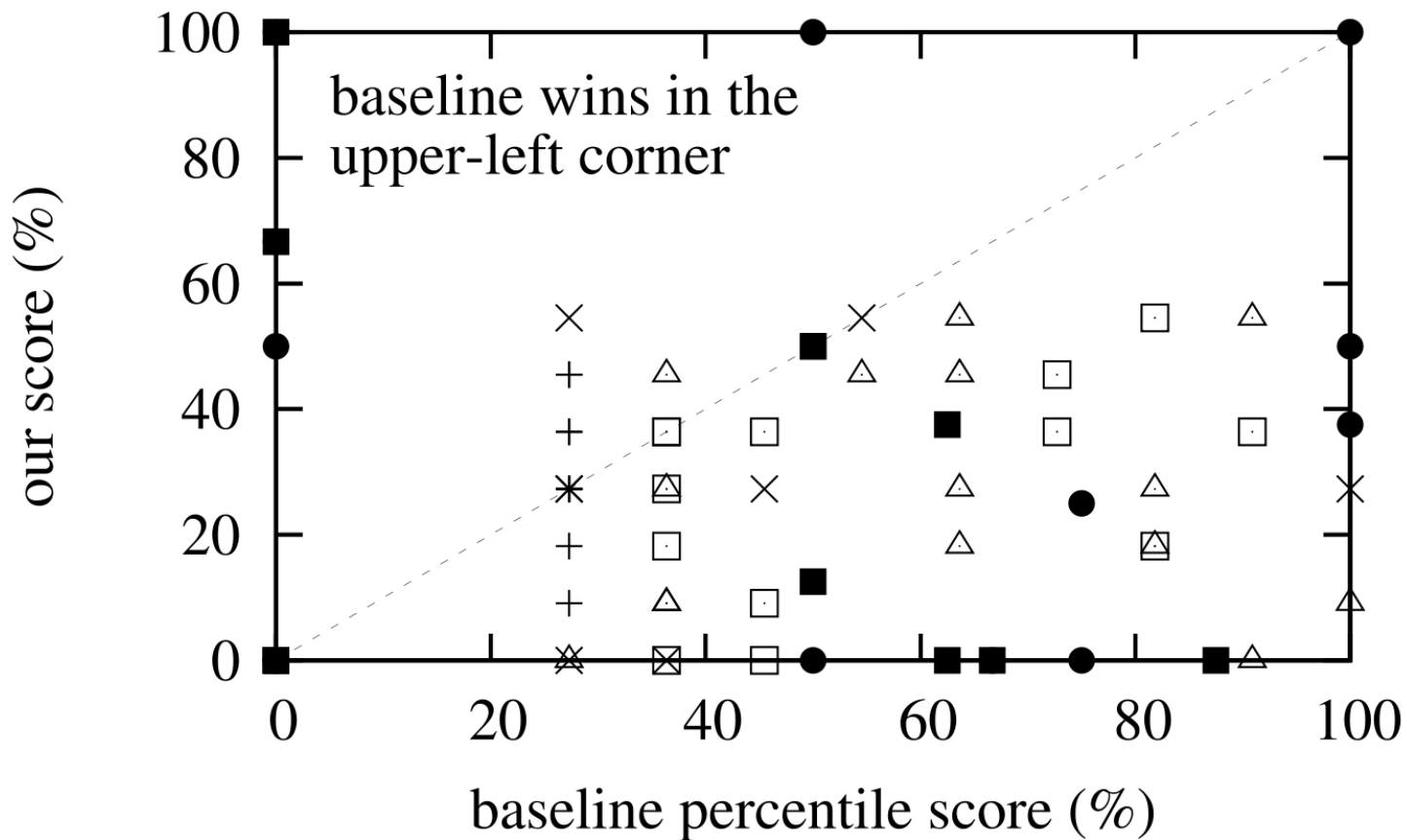


# The number of genuine causes found and missed.



Ostrowski, Mann, Sandler. LADIS '11.

# The quality of ranking relative to the K-S baseline.

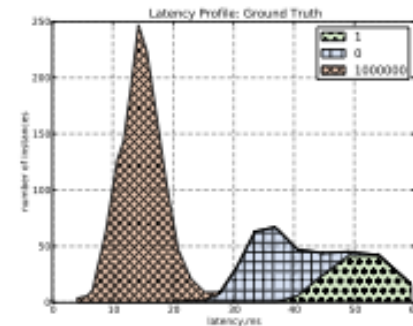
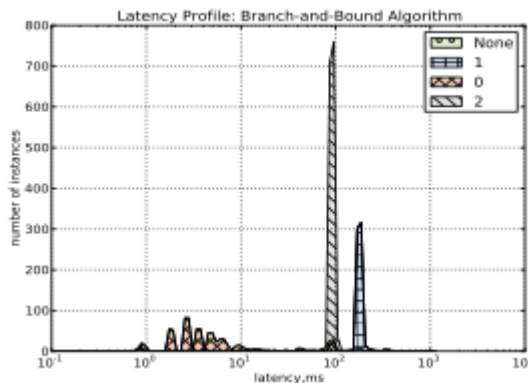
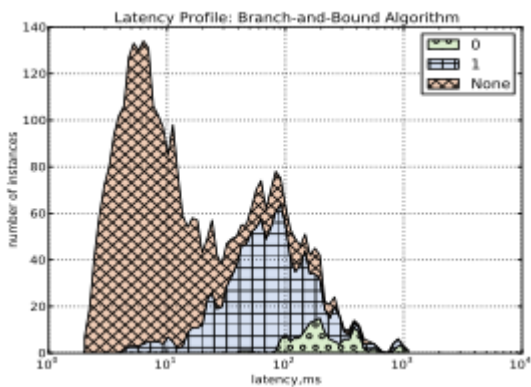
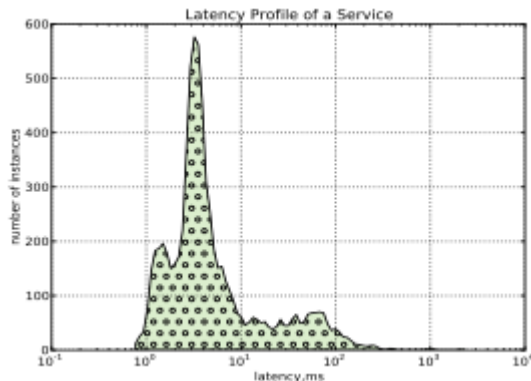


Ostrowski, Mann, Sandler. LADIS '11.

T-1/1	+	T-1/2	□	T-2/2	■
S-1/1	×	S-1/2	△	S-2/2	●

# Recent Work

## Understanding Latency Variation of Black Box Services



Krushevskaja,  
Sandler. WWW '13

# Diagnosis and Automation for Deployed Distributed Systems

Google™



- Automatic Job Health Assessment
- Real-time Diagnosis of Problem Events
- Diagnosis of Text Bugs
- Intelligent Automation