

# PRECISELY TRACE A REQUEST IN CLOUD ENVIRONMENT



Progress report meeting  
December 2013

Phuong Tran Gia  
[gia-phuong.tran@polymtl.ca](mailto:gia-phuong.tran@polymtl.ca)

Under the supervision of Prof. Michel R. Dagenais  
Dorsal Laboratory, Polytechnique Montréal

# Outline

- Research objectives
- Literature review
- Methodology
- Future work
- Challenges

# Research objectives

# Research objectives

- Research questions
- Possible contributions

# Research questions

- “How to trace and analyze a request in cloud environment at the kernel level? If impossible at the kernel level, how to minimally instrument the application’s source code and middleware?”
- “Is it possible to predict the associated overhead of tracing tools based on the traffic workloads?”

# Possible contributions

- The overall objective:
  - “Provide efficient analysis algorithms and tools in order to know exactly how a request is serviced in modern cloud-based systems and to collect the time spent on many services in the cloud”.
- Possible contributions:
  - An automated method of following a request without instrumenting the application’s source code or special RPC libraries itself in cloud environment.
  - A method to select and study user requests of interest based on the system model and other parameter, such as: set of activated probes, sampling rate or traffic workload.
  - Algorithms and views to analyze the Asynchronous RPC transaction data in cloud environment

# Literature review

# Literature review

- **Approaches to white-box tracing**
  - require the knowledge of application or middleware
  - Tracepoints
  - Request ID propagation
- **Statistical approaches to black-box tracing**
  - Events are collected from communications: incoming, outgoing messages.
  - Imprecise method.
- **Precise approaches to black-box tracing**
  - Does not require to modify the application's source code
  - Insufficient to obtain the traditional trace data.
  - 3 solutions so far: vPath<sup>[4]</sup>, BorderPatrol<sup>[5]</sup>, PreciseTracer<sup>[6]</sup>



# Literature review

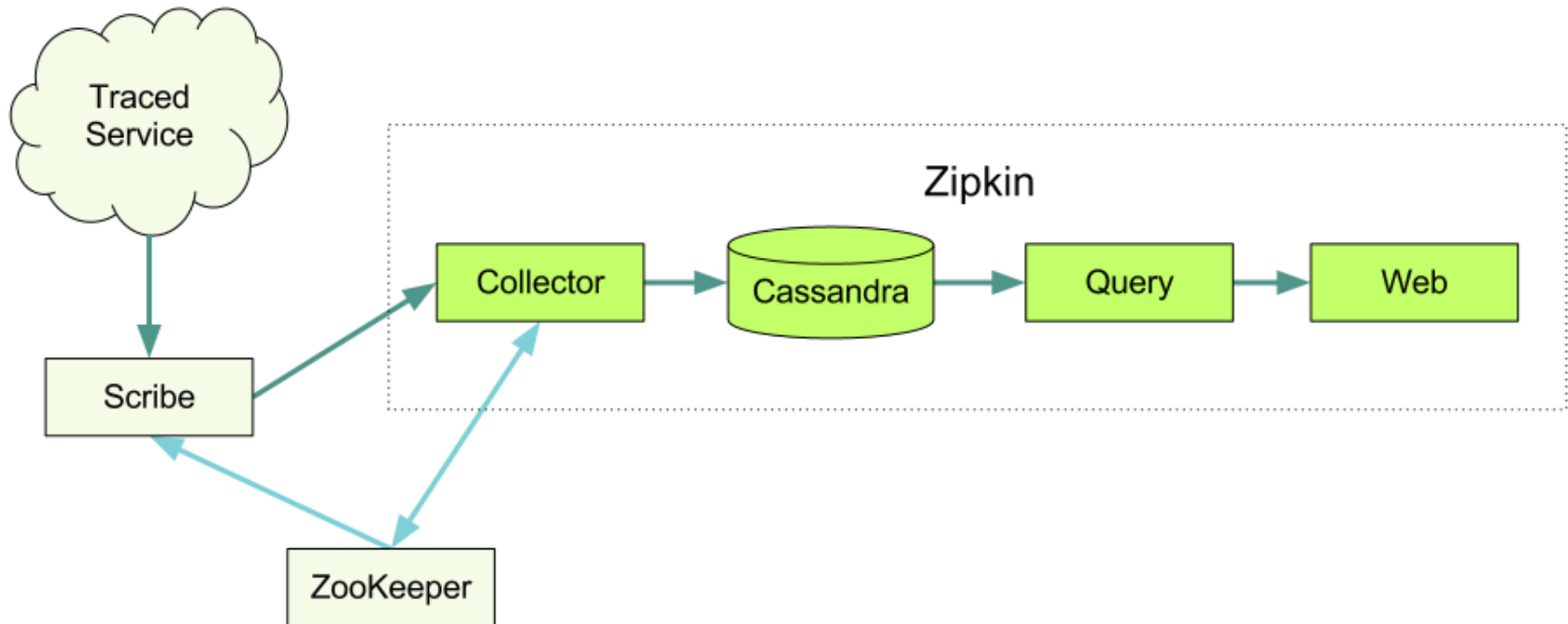
Tracing Infrastructures	Use cases	Accuracy Tracing	Sampling technique	Scalability	Operation mode
Dapper	Distributed profiling Diagnosing bottlenecks	Precise	Yes (Adaptive sampling)	Yes	Offline Online(limited)
Zipkin	Distributed profiling Diagnosing bottlenecks	Precise	Yes (Uniform sampling)	Yes	Offline
Magpie	Anomaly detection Workload modelling	Precise	No	No	Online/Offline
X-Trace	Diagnosing bottlenecks	Precise	No		
Whodunit	Distributed profiling	Precise	No		
Pip	Diagnosing bottlenecks	Precise	No	No	Offline
Pinpoint	Anomaly detection Diagnosing bottlenecks	Precise	No		Online
Stardust	Diagnosing bottlenecks Workload modelling	Precise	No		
E2Eprof	Diagnosing bottlenecks Anomaly detection	Imprecise	No	No	Online
Project 5	Diagnosing bottlenecks	Imprecise	No	No	Offline
WAP5	Diagnosing bottlenecks	Imprecise But per process	No	No	Offline
BorderPatrol	Diagnosing bottlenecks	Precise	No	No	Offline
vPath	Diagnosing bottlenecks	Precise	No	No	Offline
PreciseTracer	Diagnosing bottlenecks	Precise	Yes	Yes	Online

# Zipkin<sup>[2]</sup> – A distributed tracing framework

- Why Zipkin?
  - Helping developers gain deeper knowledge about how certain requests perform in a distributed system.
  - Helping developers gather timing data for the disparate services at Twitter.

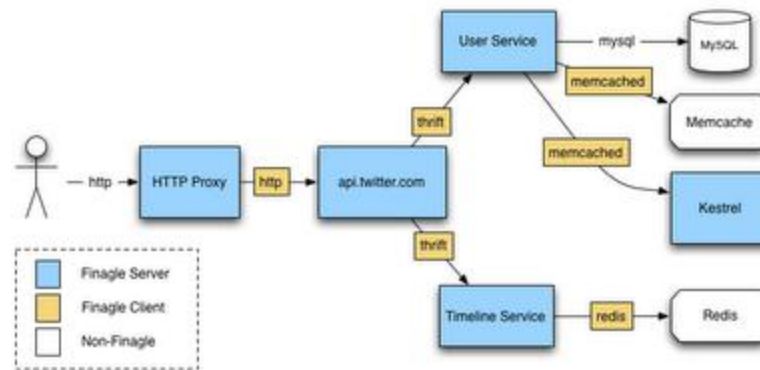


# Zipkin – Architecture



# Finagle in Zipkin

- A core module of Zipkin.
- Finagle is an asynchronous network stack for the JVM that we can use to build asynchronous Remote Procedure Call (RPC) clients and servers in Java, Scala, or any JVM-hosted language.



[github.com/twitter/finagle](https://github.com/twitter/finagle)



# Google's Dapper



dapper-2010.pdf

# Google's Dapper<sup>[1]</sup>

- Collect traces from production requests
- Low overhead
- Minimum of extra work for developer

1

**Job Selection** ?

Start Date: 05/06/2008  
 Start Hour: 09 ?  
 End Date: 05/06/2008  
 End Hour: 10 ?  
 Cluster: clusterABC  
 User: user123  
 Job: jobXYZ

**Node Information** ?

User  
 RPC or Span Name  
 Job  
 Cluster

**Cost Metric** ?

Latency ?  
 Parent Latency ?  
 Request Size ?  
 Response Size ?  
 Recursive Size ?  
 Recursive Queue Time ?

Id ?	Calls ?	Total (ms) ?	Global 90%ile Contribution (count) ?	Local 90%ile (ms) ?	Absolute Histogram (ms) ?	Scaled Histogram (ms) ?	View ?
All ?	40,990,720 (100.00%)	139,773,132.8 (100.00%)	4,098,118 (100.00%)	8.91			<a href="#">View</a>
E	3,450,880 (8.42%)	39,437,312.0 (28.22%)	1,918,437 (46.81%)	19.17			<a href="#">View</a>
R	1,658,880 (4.05%)	55,939,686.4 (40.02%)	1,658,880 (40.48%)	47.21			<a href="#">View</a>

2

**Simplified Call Tree** ?

```

  graph TD
    frontend --> search
    search --> getdocs
    getdocs --> thing1
    getdocs --> thing2
    getdocs --> helper1
    helper1 --> helper2
  
```

**Viewing Execution Pattern: E**

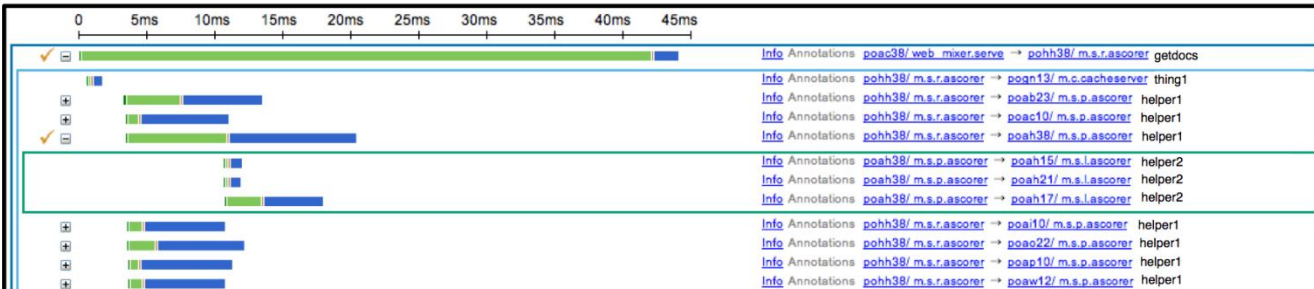
**3** **4**

Min (ms)	Max (ms)	Trace ?
51.50	66.00	<a href="#">Example</a>
40.20	51.50	<a href="#">Example</a>
31.40	40.20	<a href="#">Example</a>
24.50	31.40	<a href="#">Example</a>
24.50	31.40	<a href="#">Example</a>
24.50	31.40	<a href="#">Example</a>
24.50	31.40	<a href="#">Example</a>
24.50	31.40	<a href="#">Example</a>
19.10	24.50	<a href="#">Example</a>

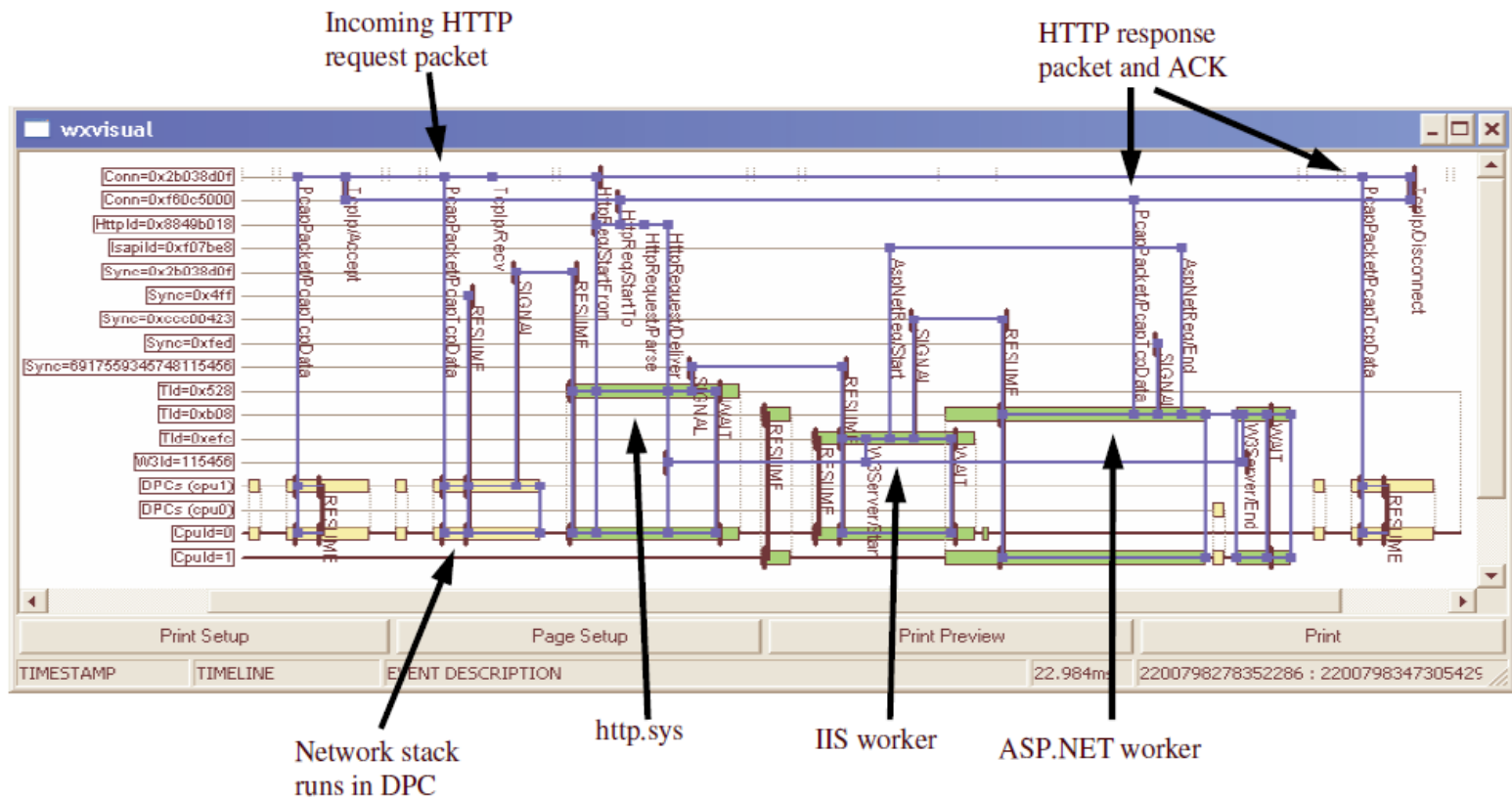


dapper-2010.pdf

5



# Magpie<sup>[3]</sup>





# Methodology

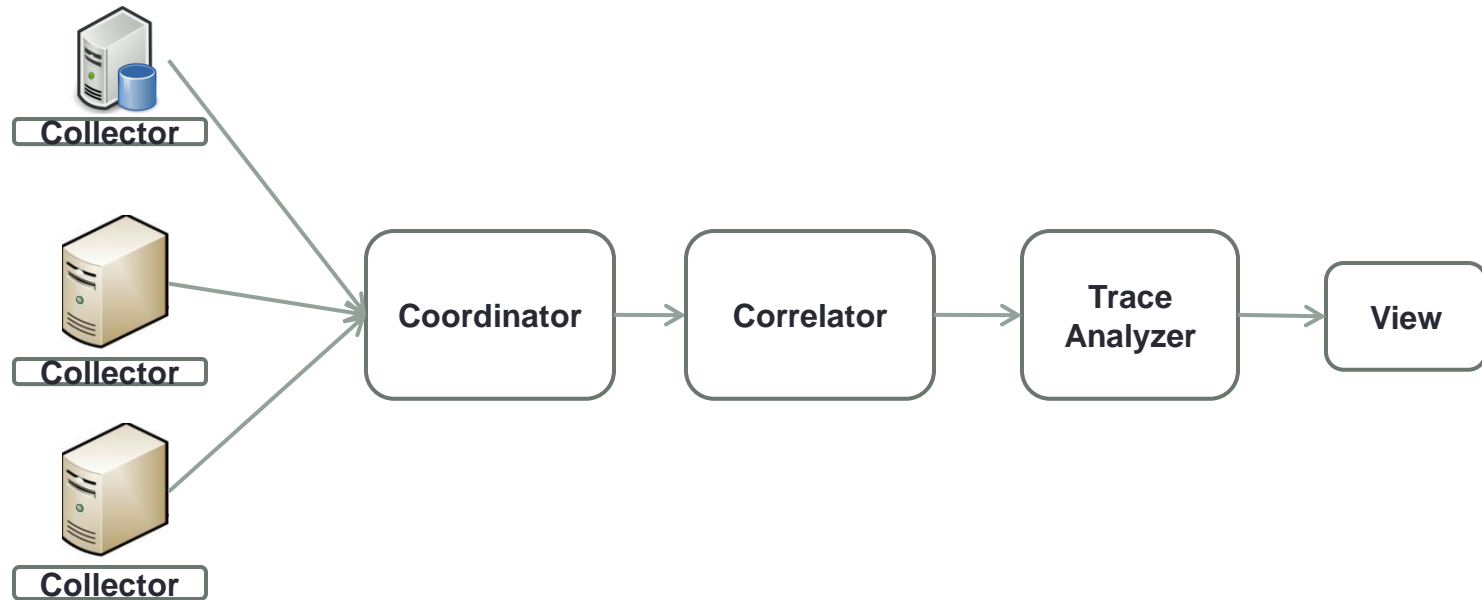
# Methodology

- System model and problem statement
- System architecture
- Temporal Join Algorithm
- Case study

# System model and problem statement

- **System model:**
  - Large clusters: disparate services, each service is running on one node.
  - Each node (component) can serve many requests in a certain period.
  - Treat each component as a grey-box (minimize the instrumentation of application source code and middleware).
- **Problem statement:**
  - Asynchronous communication between nodes (components).
  - Support event-driven architecture (like Finagle at Twitter)
  - On each component:
    - More than one threads to serve a request.
    - threads can be used many times (reused).

# System architecture



# System architecture

- Collector:
  - What information do we capture?
    - Inter-process communication between components:
      - *Domain name, Timestamp, Sender\_IP, Sender\_Port, Receiver\_IP, Receiver\_Port, Message Size*
    - Track the thread (**Tid**) performing the system call over the incoming, outgoing TCP connections.
    - Resource accounting for above thread (**Tid**) such as context switching and blocking on network.

# Approach

- How do we trace a request?
  - Design goals:
    - Determine which events from the stream pertain to a specific request?
    - Minimize the instrumentation for several types of RPC.
  - Problems:
    - The thread may post the same event for any number of different requests.
    - More than one threads to serve a request.
  - Solution: Using event-schema to describe event relationships for the particular application service, middleware of interest. (do not require to instrument the application or middleware).

# Algorithm

- Event-schema:

- Specify which attribute of event connect to other events.
- Identify each execution interval: BIND\_START, BIND\_STOP, BIND\_BASIC<sup>[3]</sup>
- Example:

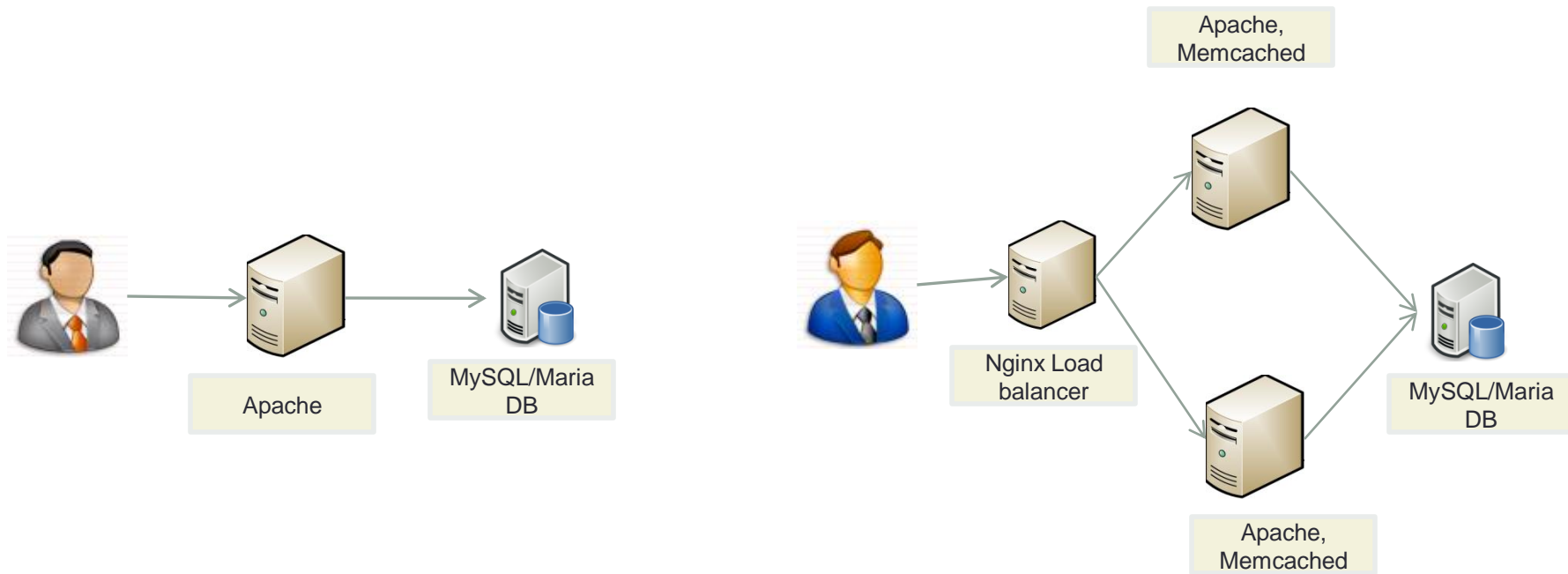
- EVENT("HttpRequest", "Start");
- ATTRIBUTE("Tid", BIND\_START, 0);
- EVENT("HttpRequest", "Start");
- ATTRIBUTE("Tid", BIND\_BASIC, 0);
- ATTRIBUTE(" **WebserverID** ", BIND\_BASIC, 0);
- EVENT(" HttpRequest ", "End");
- ATTRIBUTE("ThreadId", BIND\_BASIC, 0);
- ATTRIBUTE("**WebserverID**", BIND\_BASIC, 0);
- EVENT("HttpRequest", "Start");
- ATTRIBUTE("Tid", BIND\_STOP, 0);

- Temporal join:

- Following of a request as queries against a temporal database: each table holds the events of a given type.
- During a valid-interval (BIND\_START and BIND\_STOP) , events are joined together based on the join attribute between them (**WebserverID**) in above example of event schema.

# Case study

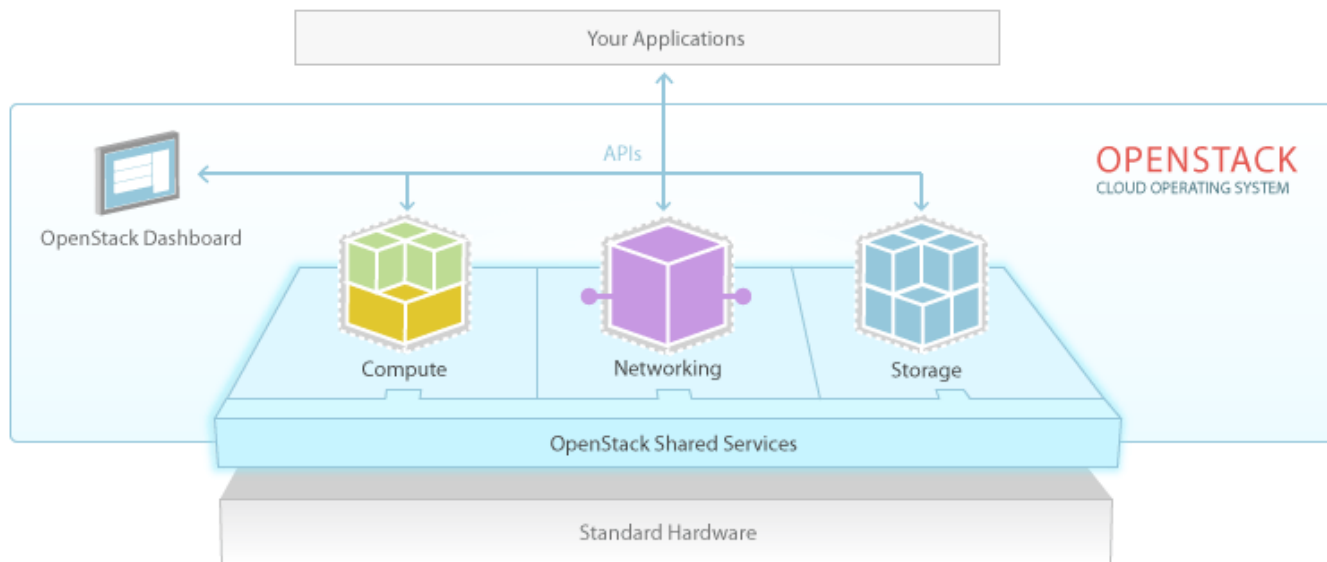
- Simple multi-tiers environment:
  - web application (mediawiki), web servers (apache), Database (MariaDB). And including Memcached.





# Case study

- Cloud platform: Openstack



# Future work

# Future work

- Build a new add-on for LTTng to trace incoming message and outgoing message at the kernel level.
- Track an incoming message within one component.
- Applying proposed approach to generate the execution path of request in a simple web application (web server farm included).
- Build a new plug-in for eclipse to see the request path.
- Implement the approach for the Openstack platform.
- Applying mathematical models (Kalman filter,...) to build a new adaptive sampling mechanism based on the traffic workloads.

# Challenges

# Challenges

- Support Synchronous/Asynchronous communication between components and within one component: a request can involve execution across threads and process boundaries within one component.
- Collect the events of Linux kernel network stack on one component.
- Simplify the event-schema as compared to Magpie.
- Improve and optimize the solutions for existing problems.

# References

- **[1]** Benjamin H. Sigelman, Luiz Andre Barroso, Mike Burrows, Pat Stephenson, Manoj Plakal, Donald Beaver, Saul Jaspán, Chandan Shanbhag, “Dapper, a Large-Scale Distributed Systems Tracing Infrastructure”, Google Technical Report dapper-2010-1, April 2010
- **[2]** Twitter Zipkin. A Distributed Tracing System, <https://github.com/twitter/zipkin>, 2012.
- **[3]** P. Barham, R. Isaacs, R. Mortier, and D. Narayanan. Magpie: on-line modelling and performance-aware systems. In 9th Workshop on Hot Topics in Operating Systems (HotOS IX), pages 85–90, May 2003.
- **[4]** B. C. Tak, C. Tang, C. Zhang, S. Govindan, B. Urgaonkar, and R. N. Chang, “vPath: precise discovery of request processing paths from blackbox observations of thread and network activities,” in Proceedings of the 2009 conference on USENIX Annual technical conference, 2009, p. 19–19.
- **[5]** E. Koskinen and J. Jannotti, “Borderpatrol: isolating events for black-box tracing,” in ACM SIGOPS Operating Systems Review, 2008, vol. 42, pp. 191–203.
- **[6]** B. Sang, J. Zhan, G. Lu, H. Wang, D. Xu, L. Wang, and Z. Zhang, “Precise, scalable, and online request tracing for multi-tier services of black boxes,” IEEE Transactions on Parallel and Distributed Systems, no. 99, p. 1–1, 2010.

**Thank you!**

**Questions?**