#### **CLOUD 2013**



The Chinese University of Hong Kong

Presented by

Zibin Zheng zbzheng@cse.cuhk.edu.hk DR<sup>2</sup>: Dynamic Request Routing for Tolerating Latency Variability in Cloud Applications

Jieming Zhu, Zibin Zheng, and Michael R. Lyu June 28, 2013



#### Introduction

- Main Challenges
- DR<sup>2</sup> Approach
- **Experiments**
- Conclusion



### Cloud computing

- Internet-based virtual computing environment
- Shared configurable resources: infrastructure, platform, software, etc.
- Pay-per-use, cost-effective

## Online cloud applications

- Search engine (e.g., Google)
- Social network (e.g., Facebook)
- E-commence website (e.g., Amazon)





### Application latency

- Time duration between a request and a response
- Evaluate the performance of online cloud applications

## The cost of latency

- 0.5s delay: 20% drop in Google's traffic
- 0.1s delay: 1% drop of Amazon's sales.



Fig. from Interxion's whitepaper



However, users perceive variability on latency, due to:

- Applications: involve numerous cloud components
- Scaling up: components deployed across data centers
- Relying on the Internet for connectivity

## Application request example



Example from Amazon's page request



# How to build consistently low-latency cloud applications, with

- Geo-distributed cloud components
- Varying latency between components

## Our proposal : Dynamic Request Routing (DR<sup>2</sup>)

- Take advantage of redundant components
  - E.g., much redundancy for fault tolerance / load balancing





#### □ A prototype of DR<sup>2</sup>



Dynamically rout the requests to different components with timely latency minimization



## Challenges



## Challenges

## Latency variability

- Relying on the Internet for connectivity
- Fluctuations over time

## Adaptivity

Adaptive to the latency dynamics

## User centricity

Optimize the request for each single user

## Scalability

Scalable and efficient



(a) Latency v.s. Time Slice





## **DR<sup>2</sup> Approach**



### DR<sup>2</sup>: Dynamic Request Routing Framework

- Phase 1: Online latency prediction
- Phase 2: Adaptive component selection



The framework of DR<sup>2</sup>



## Phase 1: Online latency prediction

#### Matrix factorization model:

|                  | $S_1$ | $S_{2}$ | $S_{3}$ | $S_4$ |                  | $S_1$ | $S_{2}$ | $S_{_3}$ | $S_4$ |
|------------------|-------|---------|---------|-------|------------------|-------|---------|----------|-------|
| $U_{1}$          | 0.2   | ?       | 0.3     | ?     | $S_1$            | 0     | ?       | 0.7      | ?     |
| $U_{_2}$         | ?     | 0.3     | ?       | 0.4   | $S_{2}$          | ?     | 0       | ?        | 0.8   |
| $U_{\rm 3}$      | 0.4   | 0.6     | ?       | ?     | $S_3$            | 0.5   | 0.6     | 0        | ?     |
| $U_{\rm 4}$      | ?     | ?       | 0.8     | 0.6   | $S_4$            | ?     | 0.4     | 0.3      | 0     |
| (a) $L^U$ Matrix |       |         |         |       | (b) $L^S$ Matrix |       |         |          |       |

$$\min \Psi(U, S, V) = \frac{1}{2} \sum_{i=1}^{n} \sum_{j=1}^{m} I_{ij}^{U} (L_{ij}^{U} - U_{i}'S_{j})^{2}$$

$$+ \frac{1}{2} \sum_{k=1}^{m} \sum_{h=1}^{m} I_{kh}^{S} (L_{kh}^{S} - V_{k}'S_{h})^{2}$$

$$+ \frac{\lambda_{U}}{2} \|U\|_{F}^{2} + \frac{\lambda_{S}}{2} \|S\|_{F}^{2} + \frac{\lambda_{V}}{2} \|V\|_{F}^{2}$$

$$Regularization term$$



## Phase 1: Online latency prediction

- Matrix factorization model
- Incremental updating of virtual coordinates (continously)

$$(u_{i}, s_{j}, L_{ij}^{U}) \longrightarrow \psi(U_{i}, S_{j}) = \frac{1}{2} (L_{ij}^{U} - U_{i}'S_{j})^{2} + \frac{\lambda_{u}}{2} \|U_{i}\|_{2}^{2} + \frac{\lambda_{s}}{2} \|S_{j}\|_{2}^{2},$$
  
$$(s_{k}, s_{h}, L_{kh}^{S}) \longrightarrow \psi(V_{k}, S_{h}) = \frac{1}{2} (L_{kh}^{S} - V_{k}'S_{h})^{2} + \frac{\lambda_{s}}{2} \|S_{h}\|_{2}^{2} + \frac{\lambda_{v}}{2} \|V_{k}\|_{2}^{2}.$$

Algorithm 1: Online Latency Prediction Algorithm **Input**: Latency data:  $L_{ij}^U$ ,  $L_{kh}^S$ **Output**: The virtual coordinates:  $U_i$ ,  $S_j$ ,  $V_k$ 1 Randomly initialize  $U \in \mathbb{R}^{d \times n}$ , and  $S, V \in \mathbb{R}^{d \times m}$ ; 2 repeat /\* Incremental updating \*/ Collect historical latency data: 3 if receive a latency data sample  $(u_i, s_j, L_{ij}^U)$  then 4 Update the virtual  $U_i \leftarrow U_i - \eta((U_i^T S_j - L_{ij}^U)S_j + \lambda_u U_i);$ 5 coordinates  $U_i$ ,  $S_i$  $S_j \leftarrow S_j - \eta((U_i^T S_j - L_{ij}^U)U_i + \lambda_s S_j);$ 6 else if receive a latency data sample  $(s_k, s_h, L_{kh}^S)$  then 7 Update the virtual  $S_h \leftarrow S_h - \eta((V_k^T S_h - L_{kh}^S)S_h + \lambda_v V_k);$ 8 coordinates  $S_h$ ,  $V_k$  $V_k \leftarrow V_k - \eta((V_k^T S_h - L_{kh}^S)V_k + \lambda_s S_h);$ 9 10 until converge;



## **DR<sup>2</sup>: Dynamic Request Routing**

## Phase 2: Adaptive component selection

#### Problem formulation

- Nodes: users and components
- Edges: available invocations
- Weights: predicted latencies
- Find shortest path for each user
- Straightforward point-point Dijkstra computation
  - Not efficient

#### Proposed solution

Get all the shortest paths in one traverse





## Phase 2: Adaptive component selection

#### Problem formulation

#### Shortest path computation (periodically update)

Algorithm 2: Adaptive Component Selection **Input**: Critical Path, Virtual coordinates:  $U_i$ ,  $S_j$ ,  $V_k$ **Output:** Component selection strategy 1 Construct the virtual graph VG based on the critical path; 2 Topologically sort VG to VG\_list; 3 foreach node v in VG list do /\* Initialization \*/ if  $v \in user$  then v.out  $\leftarrow U_i$ ; v.in  $\leftarrow$  none; 5 else v.out  $\leftarrow V_k$ ; v.in  $\leftarrow S_i$ ; 6 if v is in the last level of the critical path then 7 v.latency  $\leftarrow 0$ ; 8 else v.latency  $\leftarrow$  inf; 9 10 v.parent  $\leftarrow$  none; 11 foreach node v in VG list do foreach node w in adjacency of v do 12 if w.latency lg v.latency + w.out'\*v.in then 13 w.latency  $\leftarrow$  v.latency + w.out'\*v.in; 14 15 w.parent  $\leftarrow$  v; 16 foreach node  $v \in user$  do /\* Output the selection strategy v.path \*/ add v.parent to v.path; 17

Convert the original graph to the virtual graph (VG)

Fast shortest path algorithm on DAG with linear complexity of O(n+m)



## **Experiments**



#### Dataset description

- Dataset1: measured by ourselves
  - 1350×460 *L<sup>U</sup>*, 460×460 *L<sup>S</sup>*
- Dataset2: extracted from [Zhang et. al 2011]
  - 4532 users × 142 components × 64 time slices
- Dataset3: synthetic dataset





## Accuracy of online latency prediction

- Accuracy improves with the increasement of matrix density
- But the increasement diminishes





#### Performance comparison

- **Random**: randomly select the candidate component
- **Greedy-M**: select the best component at each step
- DR<sup>2</sup>: our approach

TABLE I. PERFORMANCE COMPARISON

| Methods                               | Density = $10\%$                    | Density = $20\%$  | Density = $30\%$  | Density = $40\%$                    | Density = $50\%$                    | Density = $100\%$ |
|---------------------------------------|-------------------------------------|-------------------|-------------------|-------------------------------------|-------------------------------------|-------------------|
| i i i i i i i i i i i i i i i i i i i | $ARE \pm std$                       | $ARE \pm std$     | $ARE \pm std$     | $ARE \pm std$                       | $ARE \pm std$                       | ARE               |
| Random                                | $6.444 \pm 0.088$                   | $6.446 \pm 0.047$ | $6.435 \pm 0.043$ | $6.431 \pm 0.035$                   | $6.405 \pm 0.069$                   | 6.436             |
| Greedy-M                              | $0.888 \pm 0.194$                   | $0.613 \pm 0.086$ | $0.517 \pm 0.110$ | $0.506 \pm 0.087$                   | $0.496 \pm 0.092$                   | 0.656             |
| $\mathbf{DR}^2$                       | $\textbf{0.412} \pm \textbf{0.108}$ | $0.269\pm0.045$   | $0.163\pm0.043$   | $\textbf{0.129} \pm \textbf{0.020}$ | $\textbf{0.089} \pm \textbf{0.025}$ | 0                 |

DR<sup>2</sup> performs the best, and is close to the baseline



### Performance on multiple users

- Randomly select 15 users as examples
- User-Noncentric/Greedy-M/DR<sup>2</sup>/Baseline



DR<sup>2</sup> optimize the component selection for each user



### Performance on multiple time slices

- **Static:** Not updating the component selection over time
- DR<sup>2</sup>: Our approach (adaptive component selection)
- Baseline: Using the exact latency data



DR<sup>2</sup> adapts to the dynamics and tolerates the latency variability

(a) Peformance on Multi. Time Slices



## **Performance Evaluation**

#### Impact of parameters

- Length of critical path
- Components per task
- Matrix density



(b) Impact: Length of Critical Path

User num. = 1350, Componet num. = 10, Matrix density = 30%



(c) Impact: Components per Task

User num. = 1350, Criti. Path Length = 10, Matrix density = 30%



(d) Impact: Matrix Density

User num. = 1350, Criti. Path Length = 10, Componet num. = 45



## Convergence time of DR<sup>2</sup>

- Batch-mode updating: periodically
- Online updating: continuously



DR<sup>2</sup> converges faster



## **Efficiency Analysis**

#### Scalability

- Number of users
- Length of the critical path
- Component number

#### Good scalability



(b) Running Time v.s. Num. of Users (c) Running Time v.s. Crit. Path Leng. (d) Running Time v.s. Compon. Num.

Criti. Path Length = 10, Componet num. = 500 Users num. = 10, Component num. = 100 User num. = 10, Criti. Path Length = 10



Conclusion & Future Work



## DR<sup>2</sup>: Dynamic request routing framework

- Tolerating latency variability in online cloud applications
- Extensive experimental results for evaluation
- Effective, adaptive, user-centric and scalable

#### Future Work

- Extend the current framework to consider load balancing strategy
- Collect more realistic application data for our experimentation

## Thank you!

Dataset available: http://www.wsdream.net