

Incorporating Heterogeneous Information for Mashup Discovery with Consistent Regularization

Yao Wan¹, Liang Chen², Qi Yu³, Tingting Liang¹ and Jian Wu¹

¹Zhejiang University, Hangzhou, China

²School of Computer Science & Information Technology, RMIT, Melbourne,

³Rochester Institute of Technology, Rochester, USA

{wanyao, cliang, liangtt, wujian2000}@zju.edu.cn, qi.yu@rit.edu

April 21, 2016



Outline

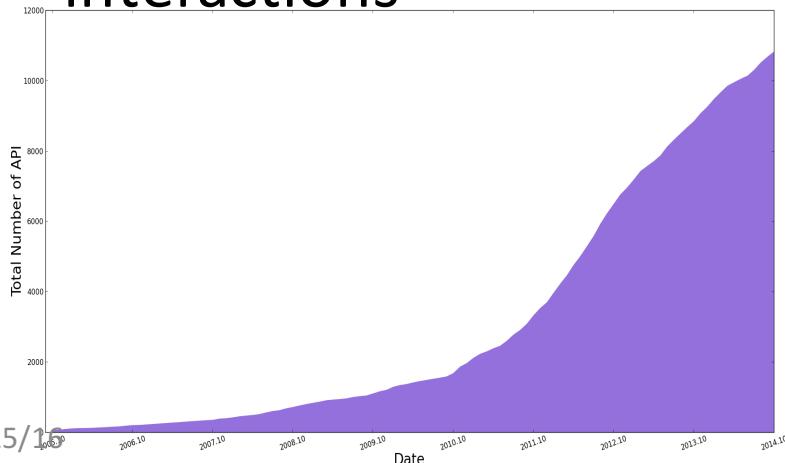
- Introduction
 - Background
 - Motivation
 - Challenge
- Model
 - A probabilistic model
 - Modeling the heterogeneous network
- Experiments
- Conclusion and Future Work



Background

□ What is web API?

- Application Programmable Interface **make web programmable**
- RESTful service
- 10,634 APIs and 6,049 mashups in ProgrammableWeb until Nov. 2014
- Increase economic transactions from web browers to API-driven interactions



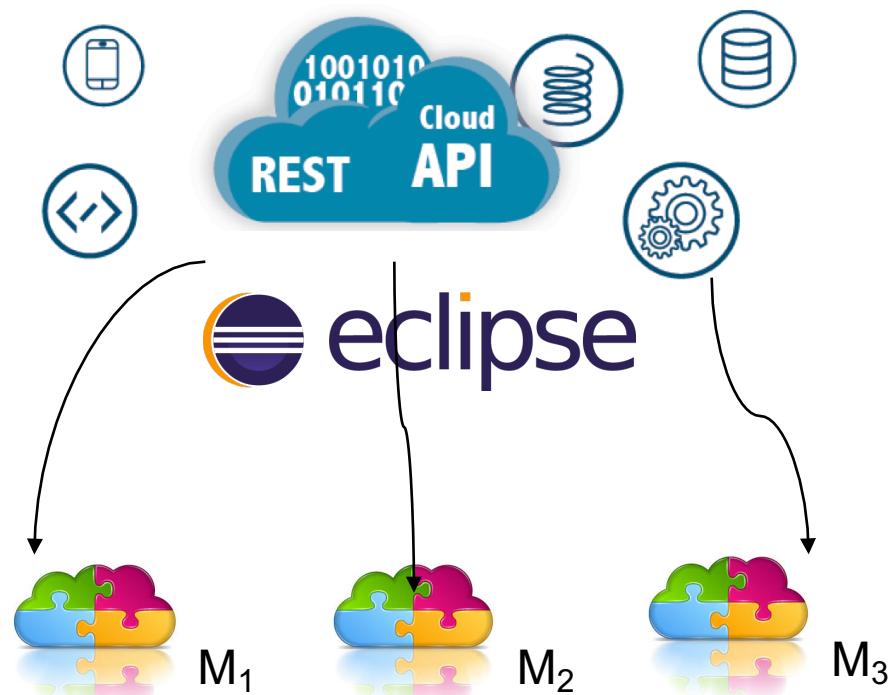


Background

□ What is Mashup? (Component Service)

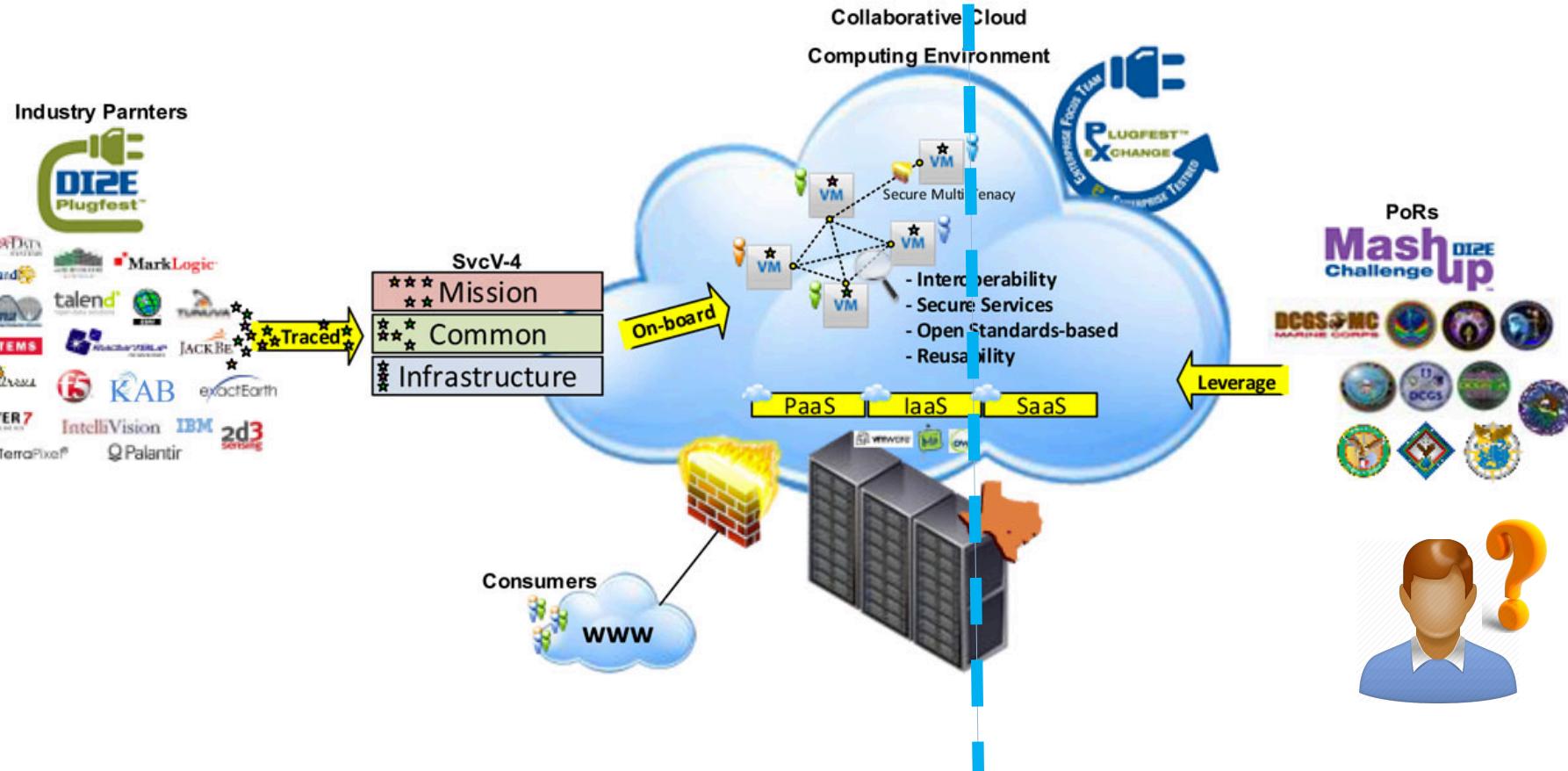
A mashup is lightweight web application created by combining capabilities from many other services

- Rapid creation (days not months)
- Reuses existing capabilities, but delivers new functions
- Request less technical skills





Motivation



Task: discover the component services - mashup



Challenges

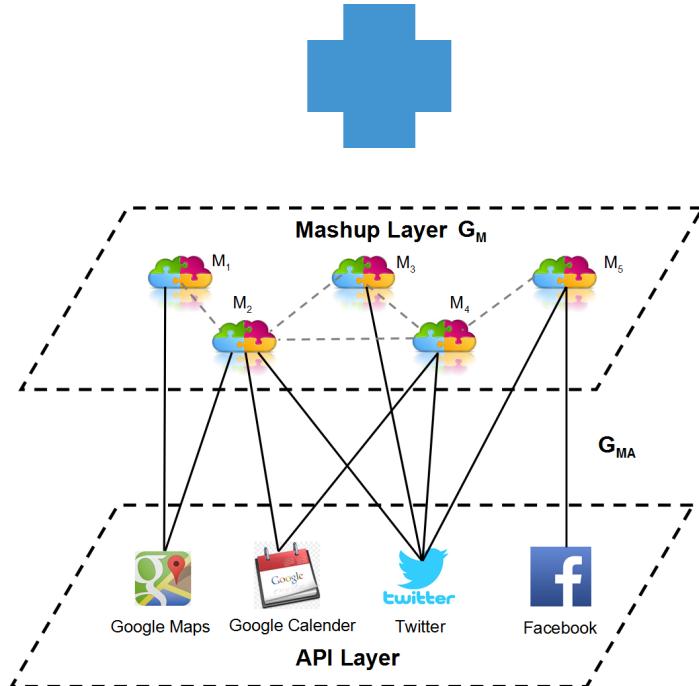
- Conventional methods
 - just utilize the semantic information of service itself
 - ignore the connection between services
- Challenge
 - Design a rank algorithm
 - Integrate the rank algorithm with the network structure

Simple Latitude Mashup

Name: Simple Latitude Mashup
Tags: Location, Family
Related APIs: Simple Latitude Open, Facebook, Twitter, Google OpenID
Description: Simple Latitude is a location-tracking iPhone app that helps users stay in touch with their families and friends...

Twitter API

Name: Twitter API
Primary category: Social
Second category: Blogging
Description: The Twitter micro-blogging service includes two RESTful APIs. The Twitter REST API methods allow developers to access core Twitter data...





A Probabilistic Model

document-centric probabilistic model: estimate the expertise of a candidate by summing the relevance of its associated documents

$$\begin{aligned} p(m_i|q) &= \lambda \sum_{a \in \mathcal{A}_{m_i}} p(m_i|a)p(a|q) + (1 - \lambda)p(m_i|q) \\ &\propto \lambda \sum_{a \in \mathcal{A}_{m_i}} p(m_i|a)p(q|a)p(a) + (1 - \lambda)p(q|m_i)p(m_i) \end{aligned}$$

- $p(m_i|a)$: the probability of API a belongs to mashup m_i
- $p(q|a)$ and $p(q|m_i)$: the semantic similarities of API and mashup
- $p(a)$ and $p(m_i)$: the quality of API a and mashup m_i

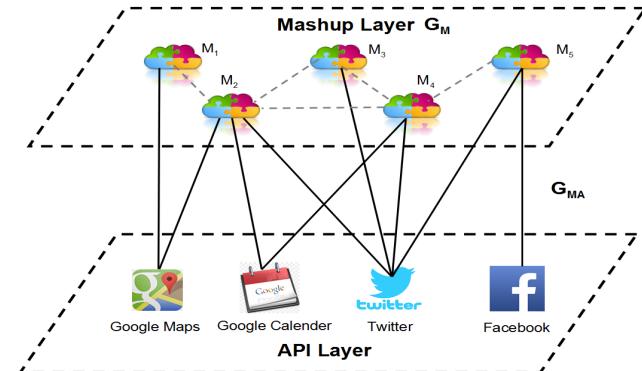
$$z = \lambda P_{MA} Q_A x + (1 - \lambda) Q_M y$$



Simple Latitude Mashup
Name: Simple Latitude Mashup
Tags: Location, Family
Related APIs: Simple Latitude Open, Facebook, Twitter, Google OpenID
Description: Simple Latitude is a location-tracking iPhone app that helps users stay in touch with their families and friends...



Twitter API
Name: Twitter API
Primary category: Social
Second category: Blogging
Description: The Twitter micro-blogging service includes two RESTful APIs. The Twitter REST API methods allow developers to access core Twitter data...





Modeling the Heterogeneous Network

□ Mashup Consistency Hypothesis:

- If two mashups share many common services with respect to a given query, then their relevance score in the queried field should be similar in some sense.

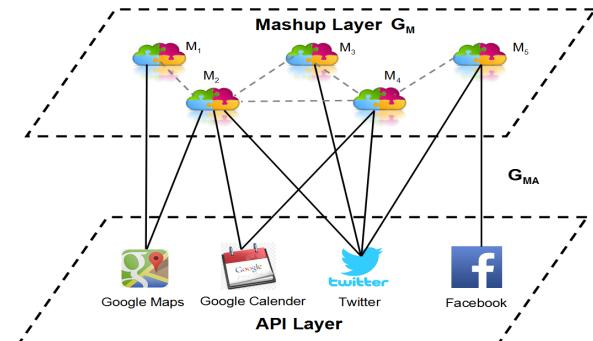
□ Objective function

$$\Omega(\mathbf{z}) = \mathbf{z}^T (\mathbf{I} - \mathbf{S}_M) \mathbf{z} + \mu \|\mathbf{z} - \mathbf{z}^0\|^2$$

$$s.t. \quad \mathbf{z}^0 = \lambda \mathbf{P}_{MA} \mathbf{Q}_A \mathbf{x} + (1 - \lambda) \mathbf{Q}_M \mathbf{y}$$

□ Similarity matrix

$$\mathbf{S}_M = \boldsymbol{\Pi}^{-1/2} \mathbf{W} \boldsymbol{\Pi}^{1/2} \quad \boldsymbol{\Pi}_{ii} = \sum_j \mathbf{W}_{ij}$$





Optimization

$$\begin{aligned}\Omega(\mathbf{z}) &= \mathbf{z}^T (\mathbf{I} - \mathbf{S}_M) \mathbf{z} + \mu \|\mathbf{z} - \mathbf{z}^0\|^2 \\ s.t. \quad \mathbf{z}^0 &= \lambda \mathbf{P}_{MA} \mathbf{Q}_A \mathbf{x} + (1 - \lambda) \mathbf{Q}_M \mathbf{y}\end{aligned}$$

- Setting $\partial\Omega(\mathbf{z})/\partial\mathbf{z} = 0$

$$(\mathbf{I} - \alpha \mathbf{S}_M) \mathbf{z}^* = (1 - \alpha) \mathbf{z}^0 \quad \alpha = 1/(1 + \mu)$$

Variant Graph Laplacian using as \mathbf{S}_M is the adjacency matrix

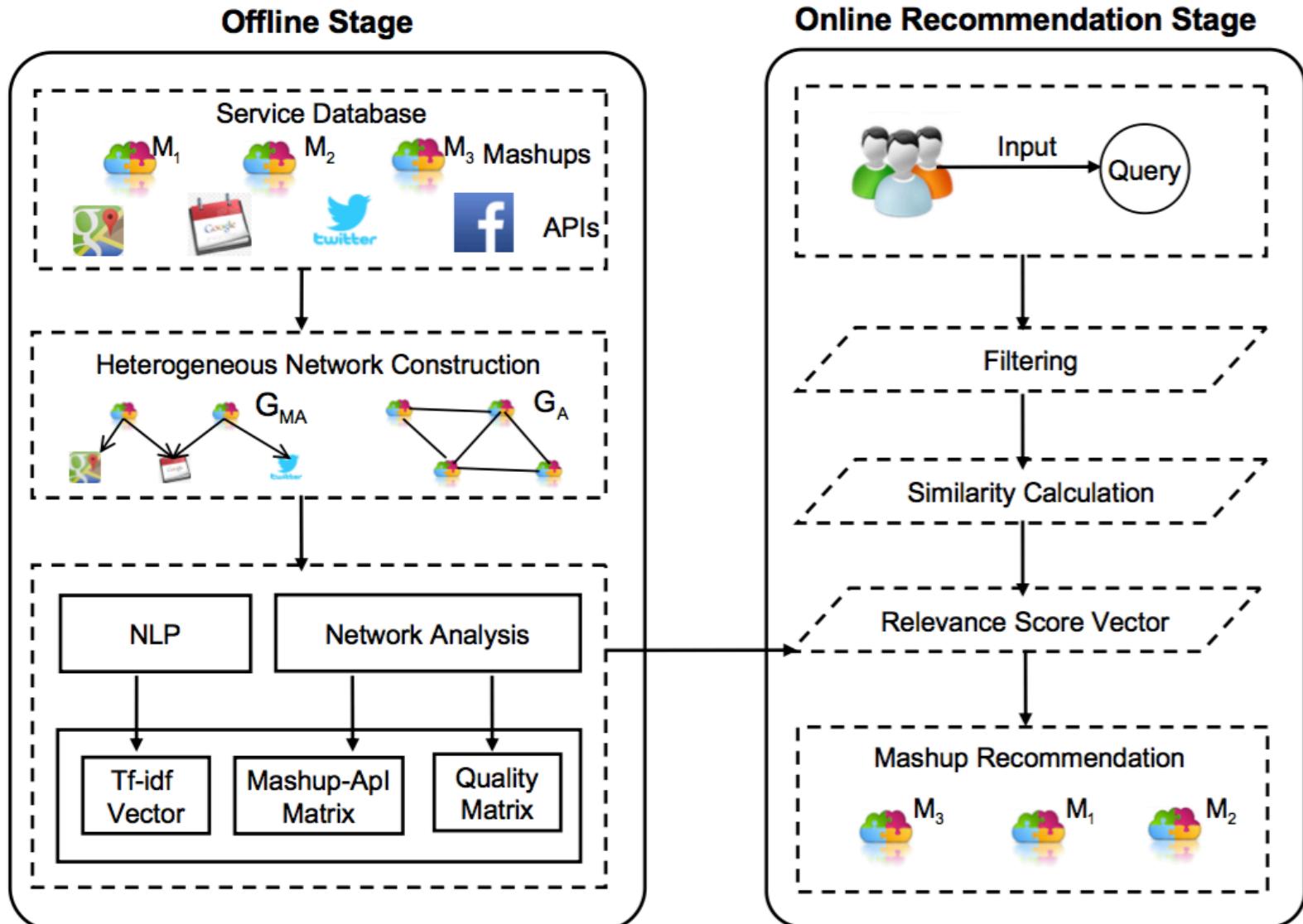
- Since \mathbf{S}_M is usually sparse

$$\mathbf{z}^{t+1} \leftarrow \beta \mathbf{S}_M \mathbf{z}^t + (1 - \beta) [\lambda \mathbf{P}_{MA} \mathbf{Q}_A \mathbf{x} + (1 - \lambda) \mathbf{Q}_M \mathbf{y}]$$

$\beta = 1/(1 + \mu)$, $\mathbf{z}^* = \mathbf{z}^\infty$ is the solution.



Implementation





Experiments

□ Dataset

Programmableweb.com



# of mashups	4,699	# edges in G_M	3,760,923
# of APIs	937	# edges in G_{MA}	8,127

□ Evaluation metrics

$$P@K = \frac{\# \text{ relevant in top } K \text{ results}}{K}$$

$$MRR = \frac{1}{|Q|} \sum_{i=1}^{|Q|} \frac{1}{rank_i}$$

$$\alpha\text{-DCG}_K = \sum_{i=1}^K \frac{G_i}{log_2(i+1)}$$



Experiments

□ Ground Truth

Category	Sample queries
Music	audio, band, lyrics, sound, music
Travel	hotel, tourism, flight, holiday, voyage, travel
Financial	bitcoin, stock, investment, price, finance, economics
eCommerce	shopping, eCommerce
Mapping	GIS, geography, GPS, mapping, traffic
Sports	NFL, fitness, sport, running, athlete
Photos	photo, image, picture, camera
Government	policy, congress, department, government, law
Game	player, game
Education	education, training, student, school
Enterprise	sales, hr, leader, enterprise
Social	media, social

The idea of binary judgment in the evaluation is: a result for a given query can be judged as **relevant** once its category is **identical with the category of the query**.



Experiments

□ Experimental results

	P@10	P@20	P@50	MRR	α -DCG
PW	0.595	0.493	0.431	-	-
MD-Sim	0.555	0.488	0.553	0.121	2.920
MD-Sim+ (vs MD-Sim)	0.575 +3.60%	0.495 +1.43%	0.534 -3.44%	0.137 +13.22%	2.916 -0.01%
MD-HIN (vs PW) (vs MD-Sim)	0.555 -6.72% 0.00%	0.53 +7.51% +8.61%	0.537 +24.59% -2.89%	0.160 - +32.23%	3.027 - +3.66%

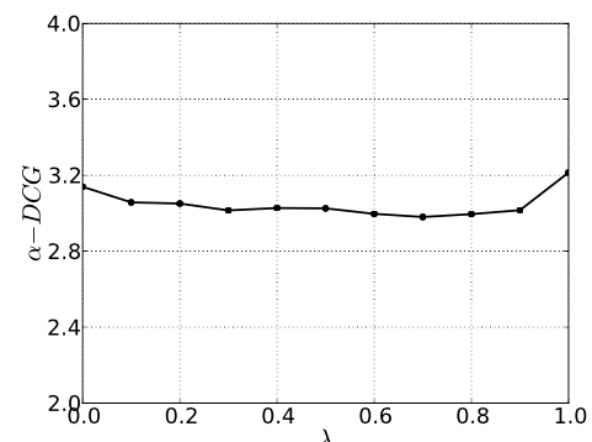
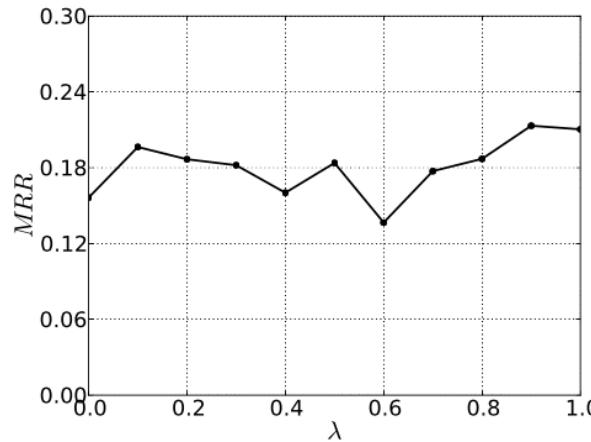
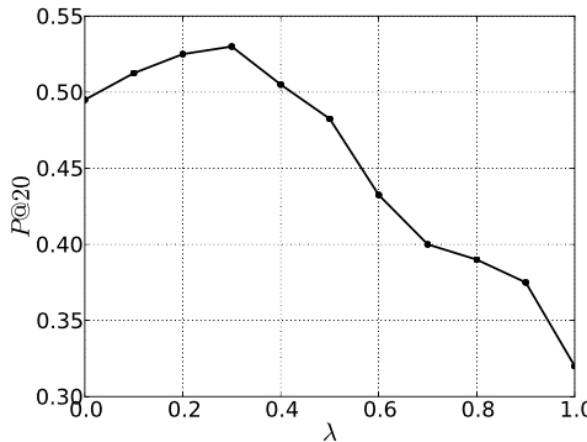
• Observations

- 1. when the value of K is large, our approach has a great advantage over ProgrammableWeb search engine.
- 2. From the perspective of ranking of results, our extended approach (MD-HIN) achieves better performance than the probabilistic approach
- 3. Among all the discovery methods, our proposed method generally achieves better performance on both metrics, indicating the effectiveness of our approach.



Parameter Analysis

□ Impact of lambda



(a) Impact of λ on $P@20$

(b) Impact of λ on MRR

(c) Impact of λ on α - DCG

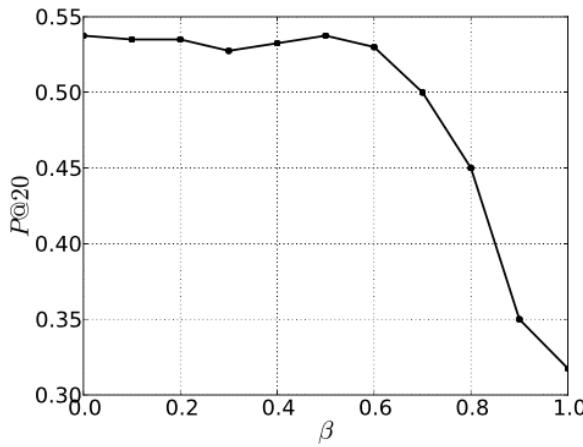
•Observations

- 1. the value of λ really has an significant impact on the performance of mashup discovery.
- 2. As λ increases, the $P@20$ value increases at first, but when λ surpasses a certain threshold, the $P@20$ value decreases with further increase of the value of λ
- 3. This phenomenon confirms the intuition that purely using the semantic information of mashups or purely employing the semantic information of their related APIs cannot generate better performance than fusing these two factors together.

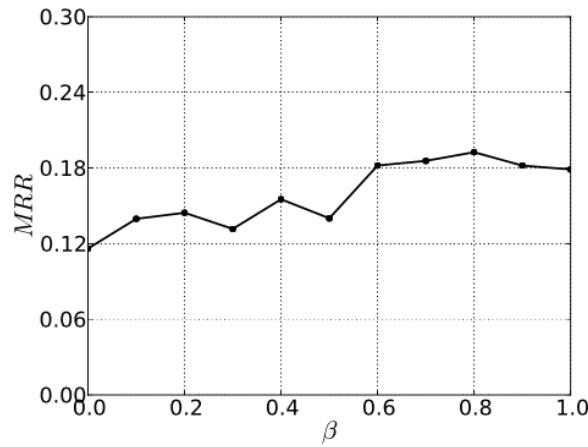


Parameter Analysis

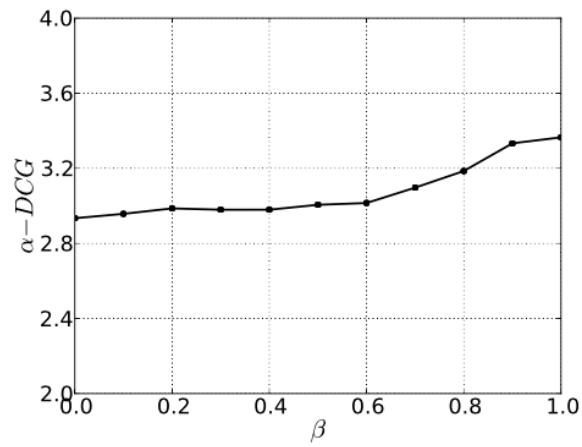
□ Impact of beta



(a) Impact of β on $P@20$



(b) Impact of β on MRR



(c) Impact of β on α - DCG

•Observations

- 1. β has a significant impact on the precision of the discovery results.



Conclusion and Future Work

□ Conclusion

- Propose an approach to improve mashup discovery by integrating the semantic information of mashups and their related APIs
- Similarity consistency is proposed and a regularization framework is implied
- Comprehensive experiments on a real-world dataset

□ Future work

- Extend our ground truth with more queries
- More evaluation metrics will be introduced

Q & A

Thank you!

wanyao@zju.edu.cn

<http://www.wanyao.me>