

# Thai Word Segmentation Web Service

The screenshot shows the THAISEMANTICS website interface. The header includes the logo and the tagline "Free Thai language resources and services". The navigation bar contains links for "หน้าหลัก", "User Profile", "Swath Word Segment", "Orchid POS Tagger", and "Web Service". The main content area is titled "หน้าหลัก > Service" and features a section for "JSON HTTP Web Service". This section explains that two methods are provided for Thai word segmentation: Swath and Part Of Speech. It also mentions that users can add custom word lists and share them. Below this text is a table titled "Support method" with two columns: "Method" and "Example Input". The table lists two methods: SWATH and ORCHID, each with a corresponding JSON example input.

Method	Example Input
SWATH	<code>{'api_key': 'YOUR API KEY', 'method': 'SWATH', 'params': ['unicode strings'], }</code>
ORCHID	<code>{'api_key': 'YOUR API KEY', 'method': 'ORCHID', 'params': [['list', 'PoS'], ['OF', 'PoS'], ['list', 'PoS'], }</code>

**Seksan Poltree (seksan.poltree@gmail.com)**  
**Asst. Prof. Kanda Saikaew (krunapon@kku.ac.th)**  
**Department of Computer Engineering**  
**Faculty of Engineering**  
**Khon Kaen University**

# Agenda

- Thai vs English text processing
- Current Thai Software and Service
- Why segmentation web service
- System Overview
- Web Application Example
- Provided Service Methods
- Comparing Service vs TLEX
- Conclusion and Future work

# Current Thai Software and Service

Resource	description	Licensing
libthai	Segmentation software + word list corpus Maximal Matching	GNU LGPL
SWATH	Segmentation software + word list corpus Maximal matching/ longest matching	GNU GPL
ORCHID	Thai Part-Of-Speech tagged corpus	NECTEC (BSD-like)
BEST	Thai segmentation solution corpus	NECTEC (BSD-like)
TLeX Service	SOAP Web service Conditional Random Field technique	Free to use

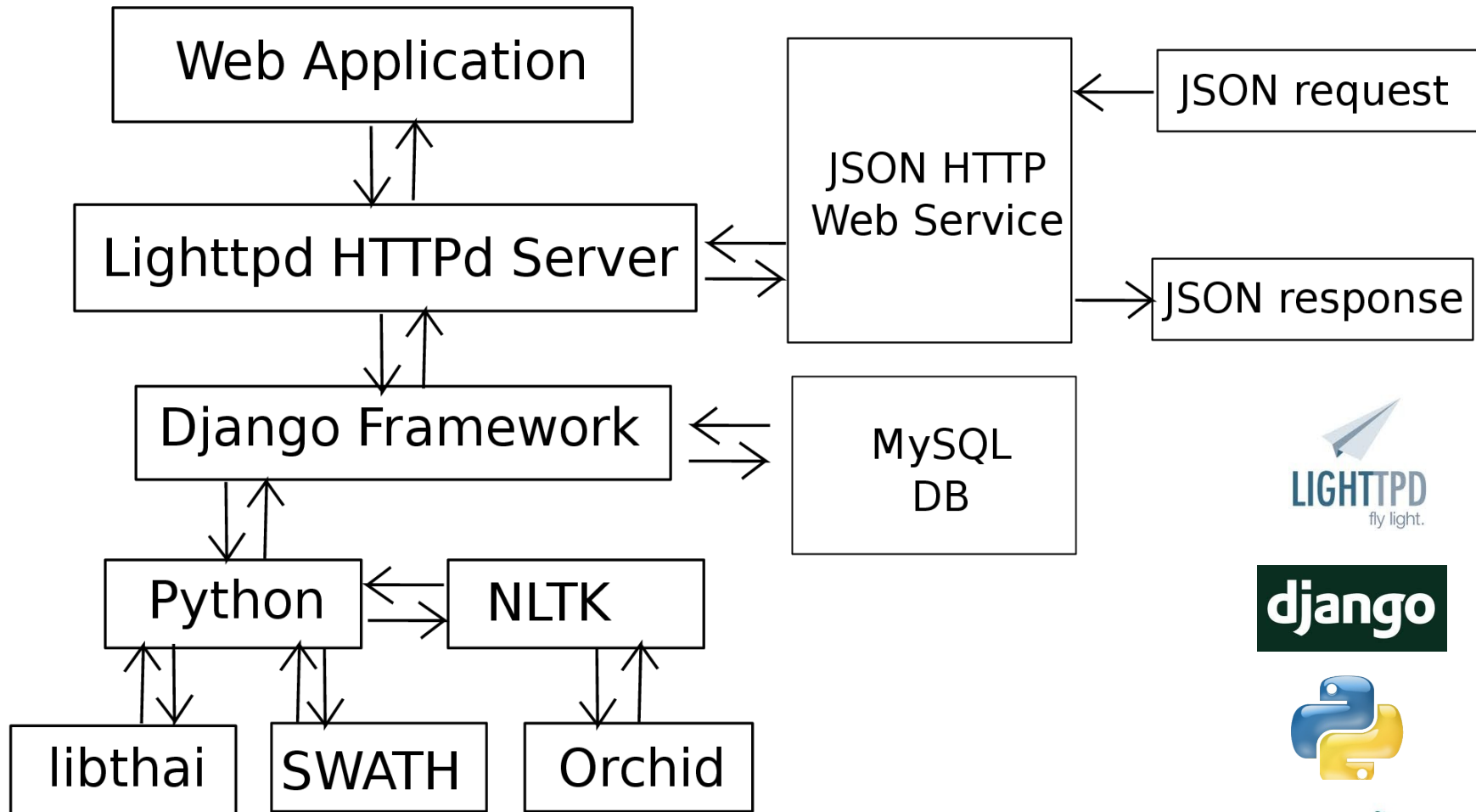


# Why Segmentation Web Service

- Increasing of web application and services
- Reducing user learning time of segmentation algorithms
- Make use of existing Thai language resources



# System Overview



# Web Application : SWATH

## Swath Word Segmentation Service

Input Sentence	Segmented Output
<p>เมื่อเปิดกระตูกกระทิงแก้มของปลาเราจะเห็นเหงือกอยู่ภายใน บนด้านหน้าของกระตูกโครงเหงือกจะมีส่วนที่ยื่นออกมาเป็นซี่เรียวยาวหรือเป็นตุ่ม ส่วนนี้เราเรียกว่า ซี่เหงือก ปลาที่กินพืชกินน้ำ เช่น กินสาหร่าย ซี่เหงือกของปลาเหล่านี้จะสั้นและมีจำนวนน้อย ส่วนปลาจำพวกที่กินแพลงก์ตอนเป็นอาหาร ซี่เหงือกยาวเรียวยาวและมีเป็นจำนวนมากสำหรับปลาบางชนิด เช่น ปลาหู ซี่เหงือกแต่ละซี่ยังแตกแขนงออกไปอีก ทั้งนี้เพื่อช่วยเพิ่มสมรรถภาพให้แก่ปลาในการกรองอาหารขนาดเล็กจากน้ำ ปลากินเนื้อหรือล่าเหยื่อเป็นอาหาร ซี่เหงือกอาจจะ</p>	<p>เมื่อ   เปิด   กระตูก   กระทิงแก้ม   ของ   ปลา   เรา   จะ   เห็น   เหงือก   อยู่   ภายใน   บน   ด้านหน้า   ของ   กระตูก   โครง   เหงือก   จะมี   ส่วน   ที่   ยื่น   ออก   มา   เป็น   ซี่   เรียวยาว   หรือ   เป็น   ตุ่ม   ส่วนนี้   เรา   เรียกว่า   ซี่   เหงือก   ปลา   ที่   กิน   พืช   กิน   น้ำ   เช่น   กิน   สาหร่าย   ซี่   เหงือก   ของ   ปลา   เหล่านี้   จะ   สั้น   และ   มี   จำนวน   น้อย   ส่วน   ปลา   จำพวก   ที่   กิน   แพลงก์   ตอน   เป็น   อาหาร   ซี่   เหงือก   ยาว   เรียวยาว   และ   มี   เป็น   จำนวน   มาก   สำหรับ   ปลา   บาง   ชนิด   เช่น   ปลาหู   ซี่   เหงือก   แต่ละ   ซี่   ยัง   แตก   แขนง   ออก   ไป  </p>
<input type="button" value="Segment"/> <input type="button" value="ล่าง"/>	

# Web Application : ORCHID

## Orchid Part of Speech Service

Input Sentence	Segmented and Tagged Output
<p>เมื่อเปิดกระตูกกระทู้แก้มของปลาเราเห็นเหงือกอยู่ ภายใน บนด้านหน้าของกระตูกโครง เหงือกจะมีส่วนที่ยื่น ออกมาเป็นซี่เรียวยาวหรือเป็นตุ่ม ส่วนนี้เราเรียกว่า ซี่ เหงือก ปลาที่กินพืชมี เช่น กินสาหร่าย ซี่เหงือกของ ปลาเหล่านี้จะสั้นและมีจำนวนน้อย ส่วนปลาจำพวกที่กิน แพลงก์ตอนเป็นอาหาร ซี่เหงือกยาวเรียวยาวและมีเป็น จำนวนมากสำหรับปลาบางชนิด เช่น ปลาหู ซี่เหงือก แต่ละซี่ยังแตกแขนงออกไปอีก ทั้งนี้เพื่อช่วยเพิ่ม สมรรถภาพให้แก่ปลาในการกรองอาหารขนาดเล็กจาก น้ำ ปลากินเนื้อหรือล่าเหยื่อเป็นอาหาร ซี่เหงือกอาจจะ</p>	<p>เมื่อ/JSBR เปิด/VACT กระตูก/NCMN กระทู้ แก้ม/UNK ของ/RPRE ปลา/NCMN เรา/PPRS  จะ/XVBM เห็น/VSTA เหงือก/UNK อยู่/XVAE  ภายใน/RPRE &lt;space&gt;/PUNC บน/RPRE ด้าน หน้า/NCMN ของ/RPRE กระตูก/NCMN โครง/UNK  เหงือก/UNK จะมี/UNK ส่วน/NCMN ที่/PREL  ยื่น/UNK ออก/XVAE มา/XVAE เป็น/VSTA ซี่/NCMN  เรียวยาว/UNK ยาว/VATT หรือ/JCRG เป็น/VSTA  ตุ่ม/UNK &lt;space&gt;/PUNC ส่วน/NCMN นี้/DDAC  เรา/PPRS เรียกว่า/VACT &lt;space&gt;/PUNC </p>

Segment and Tag

ล้าง

<http://www.thaisemantics.org/service/orchid/index>

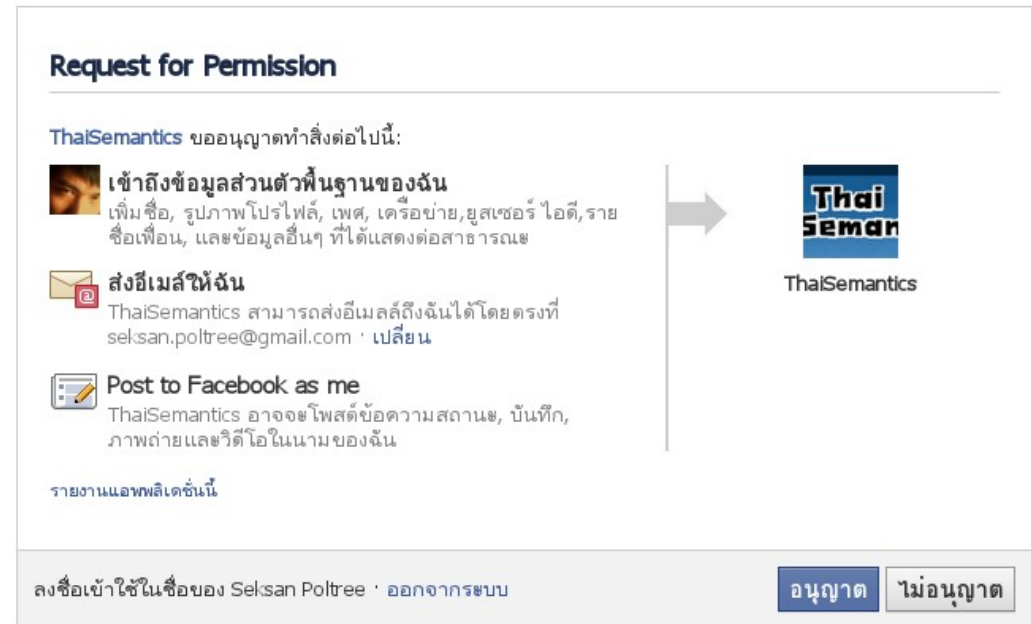


# Current Provided Service Methods

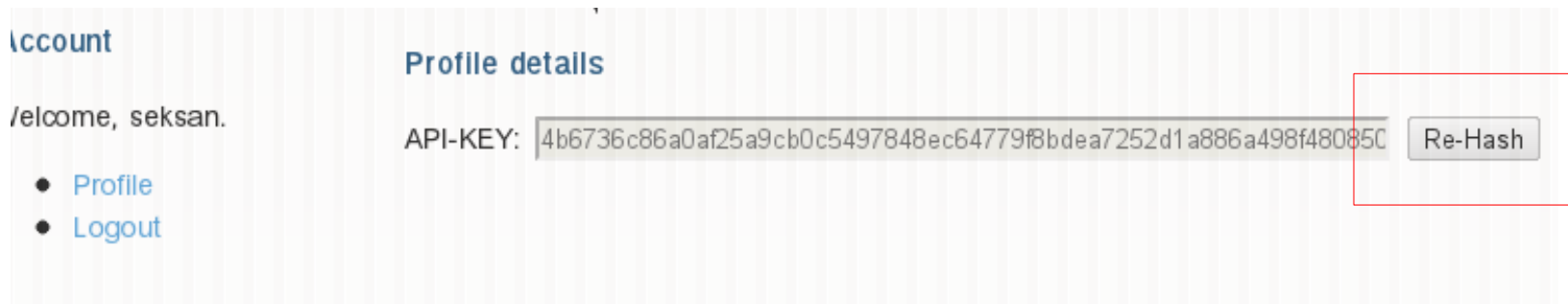
	Sample Request	Sample Response
<b>SWATH</b>	<pre>api_key': 'YOUR API KEY', 'method': 'ORCHID', 'params': [['list', 'PoS'], ['OF', 'PoS'], ['list', 'PoS'], ['list', 'PoS']], }</pre>	<pre>{"status": 0, "result": ['list', 'of', 'segmented', 'words'], }</pre>
<b>ORCHID</b>	<pre>{'api_key': 'YOUR API KEY', 'method': 'ORCHID', 'params': [['list', 'PoS'], ['OF', 'PoS'], ['list', 'PoS'], ['list', 'PoS']], }</pre>	<pre>{"status": 0, "result": [list of tagged', 'words'], }</pre>
Wrong KEY	<pre>{ 'api_key': "", 'method': 'ORCHID', 'params': ['unicode strings'], }</pre>	<pre>{"status": 1, "result": ["Wrong API key."]}</pre>
Wrong JSON	<pre>{unknown or malform json format}</pre>	<pre>{"status": -1, "result": ["Unkown request"]}</pre>

# Register to get Free API Key

- Using Facebook account instead of legacy registration



- Re-generated your API Key on demand



# Why REST, not SOAP Service?

- **REST** : REpresentational State Transfer
  - Simple, Lightweight
  - But Lack of Standard
- **SOAP** : Simple Object Access Protocol
  - XML based, Schema, Standard
  - Need more bandwidth, Higher round trip time Latency
- No complex schema description need for segmentation
- REST is more suitable!



# Why JSON not XML

- **XML** : eXtensible Markup Language
  - Self Descriptive language
  - Mark up overhead
- **JSON** : JavaScript Object Notation
  - Use simple brackets and notations
  - Suitable for simple transfer data
- No complex schema description need for segmentation , JSON is more suitable!



# Comparing Service with TLeX

## Original BEST data

เนื่องจากทั้งประเด็นเรื่อง"ความไม่เป็นธรรม"และเรื่อง"  
คู่ต่างเป็นประเด็นที่ใหญ่และซับซ้อนจึงส่งผลการสะท้อน  
แต่เพียงบางส่วนเท่านั้น |อีกทั้งประเด็นดังกล่าวมีลักษณะ  
นั้น |ข้อเสนอต่าง ๆ ที่เสนอในบทนำดังกล่าวนี้ |จึงเป็นก  
ที่ต้องการให้ท่านผู้อ่านได้นำไปขบคิดถกเถียงกันต่อไป

## SWATH Output

เนื่องจากทั้งประเด็นเรื่อง"ความไม่เป็นธรรม"และเรื่อง"  
คู่ต่างเป็นประเด็นที่ใหญ่และซับซ้อนจึงส่งผลการสะท้อน  
ได้แต่เพียงบางส่วนเท่านั้น |อีกทั้งประเด็นดังกล่าวมีลักษณะ  
ดังนั้น |ข้อเสนอต่าง ๆ ที่เสนอในบทนำดังกล่าวนี้ |จึงเป็นก  
ที่ต้องการให้ท่านผู้อ่านได้นำไปขบคิดถกเถียงกันต่อไป

## TLexs Output

เนื่องจากทั้งประเด็นเรื่อง"ความไม่เป็นธรรม"และเรื่อง"  
คู่ต่างเป็นประเด็นที่ใหญ่และซับซ้อนจึงส่งผลการสะท้อน  
ได้แต่เพียงบางส่วนเท่านั้น |อีกทั้งประเด็นดังกล่าวมีลักษณะ  
ดังนั้น |ข้อเสนอต่าง ๆ ที่เสนอในบทนำดังกล่าวนี้ |จึงเป็น  
เจตนาที่ต้องการให้ท่านผู้อ่านได้นำไปขบคิดถกเถียงกันต่อไป

- Using BEST corpora as test data
- Create simple script and call each service
- TLEX and SWATH use difference method and implementation
- Just prove of concept

# Evaluation Result

## AVERAGE SERVICE CALLING USAGE TIME

Calling Service	real	user	sys
SWATH	1m39.293s	0m5.172s	0m1.176s
ORCHID	4m50.377s	0m11.077s	0m1.212s
TLexs	3m17.045s	0m6.632s	0m1.048s

## EVALUATION RESULTS

Calling Service	Precision	Recall	F-Score
SWATH	85.93	77.62	81.57
TLexs	95.65	97.45	96.54

# Conclusion and Future work

- Create Segmentation and POS-Tagger application and services
- Create Free JSON REST Web Service
- <http://www.thaisemantics.org>
- Comparing with existing TLeX SOAP web service to prove of concept
- Include more method and corpus in the future
- Using facebook account instead of registration



# References

- [1] T. Karoonboonyanan, C. Silpa-Anan, P. Kiatisevi, P. Veerathanabutr and V. Ampornarmveth, "libthai Library". Available at: <http://linux.thai.net/projects/libthai>.
- [2] P. Charoenpornasawat, "SWATH (Smart Word Analysis for THai)". Available at: <http://www.cs.cmu.edu/~paisarn/software.html>.
- [3] T. Karoonboonyanan, "swath 0.4.1 Released". Available at: <http://linux.thai.net/svn/software/swath>.
- [4] The Royal Institute of Thailand 2525, "Thai dictionary words from the royal institute of Thailand 2525". Available at: <http://thailang.nectec.or.th/>.
- [5] V. Sornlertlamvanich, T. Charoenporn and H. Isahara, "ORCHID: Thai Part-Of-Speech Tagged Corpus". National Electronics and Computer Technology Center. Technical Report: TR-NECTEC-1997-001, 1997.
- [6] National Electronics and Computer Technology Center (NECTEC), "BEST Corpus". Available at: <http://thailang.nectec.or.th/best/>.
- [7] C. Haruechaiyasak and S. Kongyoung, "TLex: Thai Lexeme Analyser Based on the Conditional Random Fields", *Proc. 8th International Symposium on Natural Language Processing*, 2009.
- [8] National Electronics and Computer Technology Center (NECTEC), "TLex". Available at: <http://sansarn.com/tlex/>.
- [9] National Electronics and Computer Technology Center (NECTEC), "TLexs". Available at: <http://www.sansarn.com/WSeg/wsd/BnSeg.wsd/>.
- [10] K. Toutanova and C. D. Manning, "Enriching the Knowledge Sources Used in a Maximum Entropy Part-of-Speech Tagger". *Proc. the Joint SIG-DAT Conference on Empirical Methods in Natural Language Processing and Very Large Corpora*, pp. 63-70, 2000.
- [11] D. Roth and D. Zelenko, "Part of Speech Tagging Using a Network of Linear Separators", *The 17th International Conference on Computational Linguistics (1998)*, pp. 1136-1142, 1998.
- [12] The World Wide Web Consortium. "Extensible Markup Language (XML) 1.0 (Fifth Edition)". Available at : <http://www.w3.org/TR/REC-xml/>.
- [13] The World Wide Web Consortium. "SOAP Version 1.2 Part 1: Messaging Framework (Second Edition)". Available at : <http://www.w3.org/TR/soap12-part1/>.
- [14] R. T. Fielding. "Representational State Transfer (REST)". Available at : [http://www.ics.uci.edu/~fielding/pubs/dissertation/rest\\_arch\\_style.htm](http://www.ics.uci.edu/~fielding/pubs/dissertation/rest_arch_style.htm).
- [15] Ecma International, "Introducing JSON". Available at: <http://www.json.org/>
- [16] Ecma International, "Standard ECMA-262 5.1 Edition / June 2011 : ECMAScript Language Specification". Available at: <http://www.ecma-international.org/publications/files/ecma-st/ECMA-262.pdf>
- [17] Wikipedia, "Comparison of web browsers". Available at: [http://en.wikipedia.org/wiki/Comparison\\_of\\_web\\_browsers](http://en.wikipedia.org/wiki/Comparison_of_web_browsers)
- [18] C. Haruechaiyasak, S. Kongyoung and M. N. Dailey. "A Comparative Study on Thai Word Segmentation Approaches", *Proc. 5th International Conference on Electrical Engineering/Electronics, Computer, Telecommunications and Information Technology (2008)*, pp.125-128, 2008.
- [19] G. van Rossum, "Python tutorial, Technical Report CS-R9526", *Centrum voor Wiskunde en Informatica (CWI)*, Amsterdam, May, 1995.
- [20] S. Bird and E. Loper, "NLTK: The Natural Language Toolkit", *Proceedings of the ACL Workshop on Effective Tools and Methodologies for Teaching Natural Language Processing and Computational Linguistics*, pp. 62-69, 2002.
- [21] Django Software Foundation, "The Django framework". Available at : <https://www.djangoproject.com/>.
- [22] "LigHTTPD fly light". Available at : <http://www.lighttpd.net/>.
- [23] C. D. Manning, P. Raghavan and H. Schütze. "An Introduction to Information Retrieval", *Cambridge University Press, Cambridge, England*, pp.155, 2009.
- [24] G. Mulligan and D. Gračanin. "A Comparison of SOAP And REST Implementation of a Service Based Interaction Independence Middleware Framework", *Proceedings of the Winter Simulation Conference*, pp.1423-1432, 2009.
- [25] G. Wang. "Improving Data Transmission in Web Applications via the Translation between XML and JSON", *Third International Conference on Communications and Mobile Computing*, pp.182-185, 2011.





**Question?**