



Tutorial: Text Analytics for Security

Tao Xie

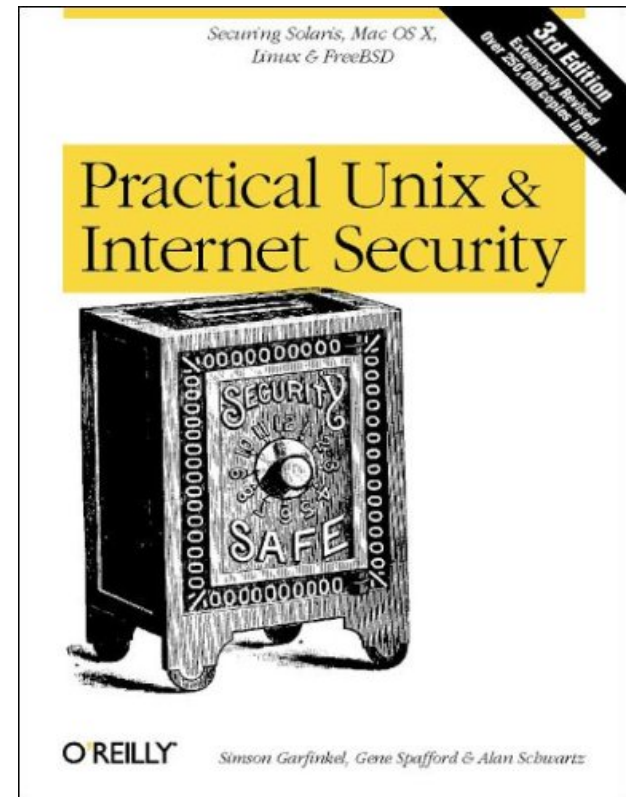
University of Illinois at Urbana-Champaign
<http://taoxie.cs.illinois.edu/>
taoxie@illinois.edu

William Enck

North Carolina State University
<http://www.enck.org>
enck@cs.ncsu.edu

What is Computer Security?

*“A computer is secure if you can depend on it and its software to behave as you **expect**.”*



User Expectations

- User expectations are a form of *context*.
- Other forms of context for security decisions
 - Temporal context (e.g., time of day)
 - Environmental context (e.g., location)
 - Execution context
 - OS level (e.g., UID, arguments)
 - Program analysis level (e.g., control flow, data flow)

Defining User Expectations

- User expectations are difficult to formally (and even informally) define.
 - Based on an individual's perception the results from past experiences and education
 - ... so, we can't be perfect
- Starting place: look at the *user interface*

Why Text Analytics?

- User interface consists of graphics and *text*
 - End users: includes finding, installing, and running the software (e.g., first run vs. subsequent)
 - Developers: includes API documentation, comments in code, and requirements documents
- Goal: *process natural language textual sources to aid security decisions*

Outline

- Introduction
- **Background on text analytics**
- Case Study 1: App Markets
- Case Study 2: ACP Rules
- Wrap-up



Challenges in Analyzing NL Data

- Unstructured
 - Hard to parse, sometimes wrong grammar
- Ambiguous: often has no defined or precise semantics (as opposed to source code)
 - Hard to understand
- Many ways to represent similar concepts
 - Hard to extract information from

```
/* We need to acquire the write IRQ lock before calling ep_unlink(). */
```

```
/* Lock must be acquired on entry to this function. */
```

```
/* Caller must hold instance lock! */
```

Why Analyzing NL Data is Easy(?)

- Redundant data
- Easy to get “good” results for simple tasks
 - Simple algorithms without much tuning effort
- Evolution/version history readily available
- Many techniques to borrow from text analytics: NLP, Machine Learning (ML), Information Retrieval (IR), etc.

Text Analytics

Knowledge Rep. &
Reasoning / Tagging

Semantic Web
Web 2.0

Search & DB

Information
Retrieval

Text Analytics

Computational
Linguistics

Natural Language
Processing

Data Analysis

Machine Learning
Text Mining

Why Analyzing NL Data is Hard(?)

- Domain specific words/phrases, and meanings
 - “Call a function” vs. call a friend
 - “Computer memory” vs. human memory
 - “This method also returns **false** if path is **null**”
- Poor quality of text
 - Inconsistent
 - grammar mistakes
 - “**true** if path is an absolute path; otherwise false” for the File class in .NET framework
 - Incomplete information

Some Major NLP/Text Analytics Tools



Text Miner

THE
POWER
TO KNOW.



BIGSHEETS



Stanford Parser

<http://nlp.stanford.edu/software/lex-parser.shtml>

SPSS
AN IBM® COMPANY

Text Analytics
for Surveys



**Unstructured
Information Management
Architecture**

An Apache Project

<http://uima.apache.org/>

<http://nlp.stanford.edu/links/statnlp.html>

<http://www.kdnuggets.com/software/text.html>

Dimensions in Text Analytics

- Three major dimensions of text analytics:
 - Representations
 - ...from words to partial/full parsing
 - Techniques
 - ...from manual work to learning
 - Tasks
 - ...from search, over (un-)supervised learning, summarization, ...

Major Text Representations

- Words (stop words, stemming)
- Part-of-speech tags

- Chunk parsing (chunking)
- Semantic role labeling
- Vector space model

Words' Properties

- Relations among word surface forms and their senses:
 - **Homonymy**: same form, but different meaning (e.g. bank: river bank, financial institution)
 - **Polysemy**: same form, related meaning (e.g. bank: blood bank, financial institution)
 - **Synonymy**: different form, same meaning (e.g. singer, vocalist)
 - **Hyponymy**: one word denotes a subclass of another (e.g. breakfast, meal)
- General thesaurus: WordNet, existing in many other languages (e.g. EuroWordNet)
 - <http://wordnet.princeton.edu/>
 - <http://www.illc.uva.nl/EuroWordNet/>

Stop Words

- Stop words are words that from non-linguistic view do not carry information
 - ...they have mainly functional role
 - ...usually we remove them to help mining techniques to perform better
- Stop words are language dependent – examples:
 - **English**: A, ABOUT, ABOVE, ACROSS, AFTER, AGAIN, AGAINST, ALL, ALMOST, ALONE, ALONG, ALREADY, ...

Stemming

- Different forms of the same word are usually problematic for text analysis, because they have **different spelling and similar meaning** (e.g. learns, learned, learning,...)
- **Stemming** is a process of transforming a word into its stem (normalized form)
 - ...stemming provides an inexpensive mechanism to merge

Stemming cont.

- For English is mostly used Porter stemmer at <http://www.tartarus.org/~martin/PorterStemmer/>
- Example cascade rules used in English Porter stemmer
 - ATIONAL -> ATE relational -> relate
 - TIONAL -> TION conditional -> condition
 - ENCI -> ENCE valenci -> valence
 - ANCI -> ANCE hesitanci -> hesitance
 - IZER -> IZE digitizer -> digitize
 - ABLI -> ABLE conformabli -> conformable
 - ALLI -> AL radicalli -> radical
 - ENTLI -> ENT differentli -> different
 - ELI -> E vileli -> vile
 - OUSLI -> OUS analogousli -> analogous

Part-of-Speech Tags

- Part-of-speech tags specify word types enabling to differentiate words functions
 - For text analysis, part-of-speech tag is used mainly for “information extraction” where we are interested in e.g., named entities (“noun phrases”)
 - Another possible use is reduction of the vocabulary (features)
 - ...it is known that nouns carry most of the information in text documents
- Part-of-Speech taggers are usually learned on manually tagged data

Part-of-Speech Table

part of speech	function or "job"	example words	example sentences
<u>Verb</u>	action or state	(to) be, have, do, like, work, sing, can, must	EnglishClub.com is a web site. I like EnglishClub.com.
<u>Noun</u>	thing or person	pen, dog, work, music, town, London, teacher, John	This is my dog . He lives in my house . We live in London .
<u>Adjective</u>	describes a noun	a/an, the, 69, some, good, big, red, well, interesting	My dog is big . I like big dogs.
<u>Adverb</u>	describes a verb, adjective or adverb	quickly, silently, well, badly, very, really	My dog eats quickly . When he is very hungry, he eats really quickly.
<u>Pronoun</u>	replaces a noun	I, you, he, she, some	Tara is Indian. She is beautiful.
<u>Preposition</u>	links a noun to another word	to, at, after, on, but	We went to school on Monday.
<u>Conjunction</u>	joins clauses or sentences or words	and, but, when	I like dogs and I like cats. I like cats and dogs. I like dogs but I don't like cats.
<u>Interjection</u>	short exclamation, sometimes inserted into a sentence	oh!, ouch!, hi!, well	Ouch! That hurts! Hi! How are you? Well, I don't know.

Part-of-Speech Examples

verb
Stop!

noun	verb
John	works.

noun	verb	verb
John	is	working.

pronoun	verb	noun
She	loves	animals.

noun	verb	adjective	noun
Animals	like	kind	people.

noun	verb	noun	adverb
Tara	speaks	English	well.

noun	verb	adjective	noun
Tara	speaks	good	English.

pronoun	verb	preposition	adjective	noun	adverb
She	ran	to	the	station	quickly.

pron.	verb	adj.	noun	conjunction	pron.	verb	pron.
She	likes	big	snakes	but	I	hate	them.

Here is a sentence that contains every part of speech:

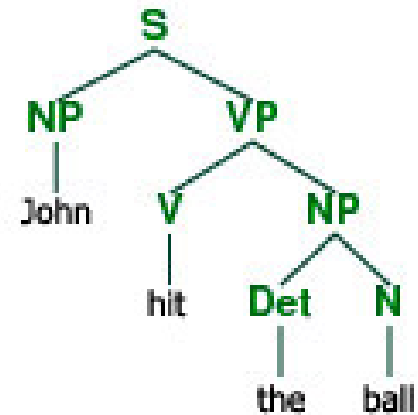
interjection	pron.	conj.	adj.	noun	verb	prep.	noun	adverb
Well,	she	and	young	John	walk	to	school	slowly.

Part of Speech Tags

Tag	Description	Example	Tag	Description	Example
CC	Coordin. Conjunction	<i>and, but, or</i>	SYM	Symbol	<i>+, %, &</i>
CD	Cardinal number	<i>one, two, three</i>	TO	"to"	<i>to</i>
DT	Determiner	<i>a, the</i>	UH	Interjection	<i>ah, oops</i>
EX	Existential 'there'	<i>there</i>	VB	Verb, base form	<i>eat</i>
FW	Foreign word	<i>mea culpa</i>	VBD	Verb, past tense	<i>ate</i>
IN	Preposition/sub-conj	<i>of, in, by</i>	VBG	Verb, gerund	<i>eating</i>
JJ	Adjective	<i>yellow</i>	VBN	Verb, past participle	<i>eaten</i>
JJR	Adj., comparative	<i>bigger</i>	VBP	Verb, non-3sg pres	<i>eat</i>
JJS	Adj., superlative	<i>wildest</i>	VBZ	Verb, 3sg pres	<i>eats</i>
LS	List item marker	<i>1, 2, One</i>	WDT	Wh-determiner	<i>which, that</i>
MD	Modal	<i>can, should</i>	WP	Wh-pronoun	<i>what, who</i>
NN	Noun, sing. or mass	<i>llama</i>	WP\$	Possessive wh-	<i>whose</i>
NNS	Noun, plural	<i>llamas</i>	WRB	Wh-adverb	<i>how, where</i>
NNP	Proper noun, singular	<i>IBM</i>	\$	Dollar sign	<i>\$</i>
NNPS	Proper noun, plural	<i>Carolinas</i>	#	Pound sign	<i>#</i>
PDT	Predeterminer	<i>all, both</i>	"	Left quote	<i>(' or ")</i>
POS	Possessive ending	<i>'s</i>	"	Right quote	<i>(' or ")</i>
PP	Personal pronoun	<i>I, you, he</i>	(Left parenthesis	<i>([({ <)</i>
PP\$	Possessive pronoun	<i>your, one's</i>)	Right parenthesis	<i>(]) } >)</i>
RB	Adverb	<i>quickly, never</i>	,	Comma	<i>,</i>
RBR	Adverb, comparative	<i>faster</i>	.	Sentence-final punc	<i>(. ! ?)</i>
RBS	Adverb, superlative	<i>fastest</i>	:	Mid-sentence punc	<i>(: ; ... - -)</i>
RP	Particle	<i>up, off</i>			

Full Parsing

- Parsing provides maximum structural information per sentence
- Input: a sentence → output: a parse tree
- For most text analysis techniques, the information in parse trees is too complex
- Problems with full parsing:
 - Low accuracy
 - Slow
 - Domain Specific





Chunk Parsing

- Break text up into non-overlapping contiguous subsets of tokens.
 - aka. partial/shallow parsing, light parsing.
- What is it useful for?
 - Entity recognition
 - people, locations, organizations
 - Studying linguistic patterns
 - gave NP
 - gave up NP in NP
 - gave NP NP
 - gave NP to NP
 - Can ignore complex structure when not relevant



Chunk Parsing

Goal: divide a sentence into a sequence of chunks.

- Chunks are non-overlapping regions of a text

[I] saw [a tall man] in [the park]

- Chunks are non-recursive
 - A chunk cannot contain other chunks
- Chunks are non-exhaustive
 - Not all words are included in the chunks

Chunk Parsing Techniques

- Chunk parsers usually ignore lexical content
- Only need to look at part-of-speech tags
- Techniques for implementing chunk parsing
 - E.g., Regular expression matching

Regular Expression Matching

- Define a regular expression that matches the sequences of tags in a chunk

– A simple noun phrase chunk regexp:

`<DT> ? <JJ> * <NN.??>`

- Chunk all matching subsequences:

The /DT little /JJ cat /NN sat /VBD on /IN the /DT mat /NN

[The /DT little /JJ cat /NN] sat /VBD on /IN [the /DT mat /NN]

- If matching subsequences overlap, the first one gets priority

DT: Determiner JJ: Adjective NN: Noun, sing, or mass

VBD: Verb, past tense IN: Preposition/sub-conj Verb



Semantic Role Labeling

Giving Semantic Labels to Phrases

- [AGENT John] **broke** [THEME the window]
- [THEME The window] **broke**
- [AGENT Sotheby's] .. **offered** [RECIPIENT the Dorrance heirs]
[THEME a money-back guarantee]
- [AGENT Sotheby's] **offered** [THEME a money-back guarantee] to
[RECIPIENT the Dorrance heirs]
- [THEME a money-back guarantee] **offered** by [AGENT Sotheby's]
- [RECIPIENT the Dorrance heirs] will [ARM-NEG not]
be **offered** [THEME a money-back guarantee]

Semantic Role Labeling Good for *Question Answering*

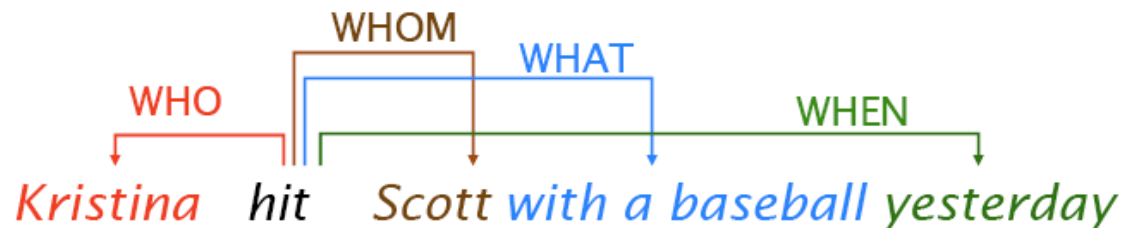
Q: What was the name of the first computer system that defeated Kasparov?

A: [PATIENT Kasparov] was defeated by [AGENT Deep Blue] [TIME in 1997].

Q: When was Napoleon defeated?

Look for: [PATIENT Napoleon] [PRED defeat-synset] [ARGM-TMP *ANS*]

More generally:



- **Who** hit Scott with a baseball?
- **Whom** did Kristina hit with a baseball?
- **What** did Kristina hit Scott with?
- **When** did Kristina hit Scott with a baseball?

Typical Semantic Roles

Thematic Role	Definition
AGENT	The volitional causer of an event
EXPERIENCER	The experiencer of an event
FORCE	The non-volitional causer of the event
THEME	The participant most directly affected by an event
RESULT	The end product of an event
CONTENT	The proposition or content of a propositional event
INSTRUMENT	An instrument used in an event
BENEFICIARY	The beneficiary of an event
SOURCE	The origin of the object of a transfer event
GOAL	The destination of an object of a transfer event

Example Semantic Roles

Thematic Role	Example
AGENT	<i>The waiter</i> spilled the soup.
EXPERIENCER	<i>John</i> has a headache.
FORCE	<i>The wind</i> blows debris from the mall into our yards.
THEME	Only after Benjamin Franklin broke <i>the ice</i> ...
RESULT	The French government has built a <i>regulation-size baseball diamond</i> ...
CONTENT	Mona asked “ <i>You met Mary Ann at a supermarket?</i> ”
INSTRUMENT	He turned to poaching catfish, stunning them <i>with a shocking device</i> ...
BENEFICIARY	Whenever Ann Callahan makes hotel reservations <i>for her boss</i> ...
SOURCE	I flew in <i>from Boston</i> .
GOAL	I drove <i>to Portland</i> .

Outline

- Introduction
- Background on text analytics
- **Case Study 1: App Markets**
- Case Study 2: ACP Rules
- Wrap-up

Case Study: App Markets

- App Markets have played an important role in the popularity of mobile devices
- Provide users with a textual description of each application's functionality



Apple App Store



Google Play



Microsoft Windows Phone

Current Practice

- Apple: **market's** responsibility
 - Apple performs manual inspection
- Google: **user's** responsibility
 - Users approve permissions for security/privacy
 - Bouncer (static/dynamic malware analysis)
- Windows Phone: hybrid
 - Permissions / manual inspection

Is Program Analysis Sufficient?

- Previous approaches look at permissions, code, and runtime behaviors
- Caveat: *what does the user expect?*
 - GPS Tracker: record and send location
 - Phone-call Recorder: record audio during call
 - One-Click Root: exploit vulnerability
 - Others are more subtle



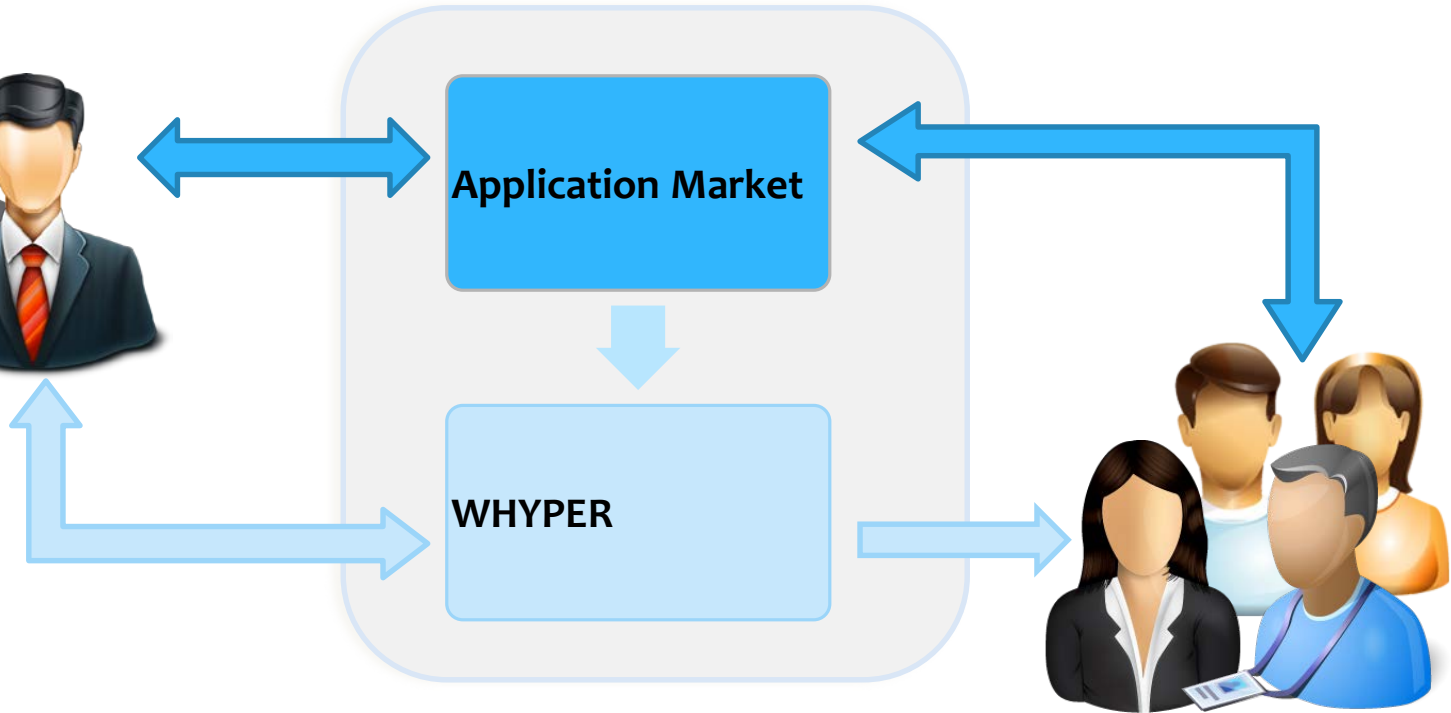
Vision

- Goal: *bridge gap between user expectation and app behavior*
- WHYPER is a first step in this direction
- Focus on permission and app descriptions
 - Limited to permissions that protect “user understandable” resources



WHYPER Overview

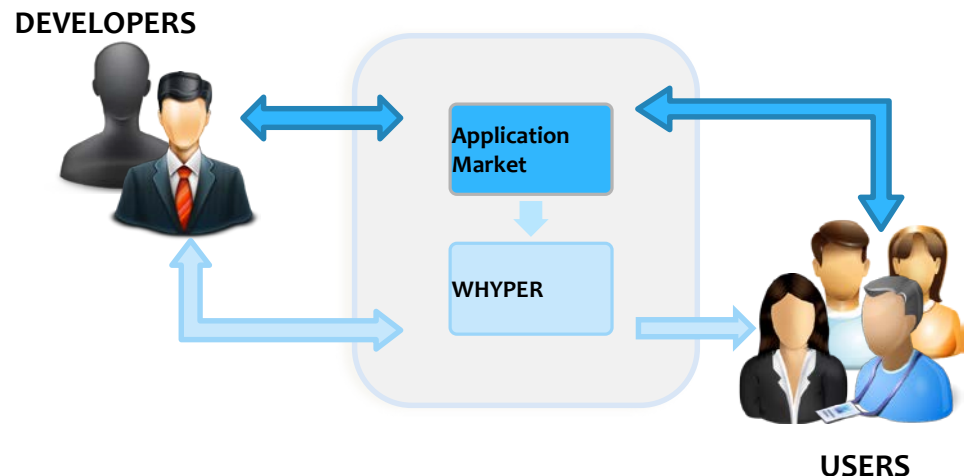
DEVELOPERS



USERS

Use Cases

- Enhance user experience while installing apps
- Functionality disclosure to during application submission to market
- Complementing program analysis to ensure more appropriate justifications



Straw man: Keyword Search

- Confounding effects:

- Certain keywords such as “contact” have a confounding meaning, e.g.,

“... displays user contacts, ...” vs “... contact me at abc@xyz.com”

- Semantic Interference:

- Sentences often describe a sensitive operation such as reading contacts without actually referring to the keyword “contact”, e.g.,

“share yoga exercises with your friends via email, sms”

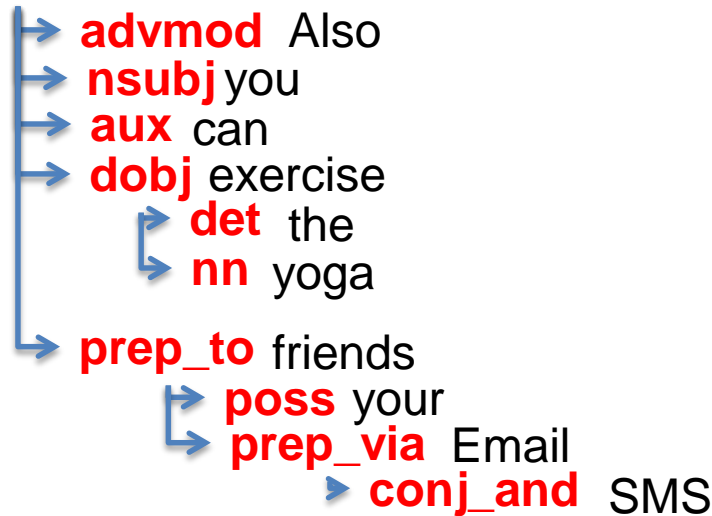
Preprocessor

- Period Handling
 - Decimals, ellipsis, shorthand notations (Mr., Dr.)
- Sentence Boundaries
 - Tabs, bullet points, delimiters (:)
 - Symbols (*,-) and enumeration sentence
- Named Entity Handling
 - E.g., “Pandora internet radio”
- Abbreviation Handling
 - E.g., “Instant Message (IM)”

Intermediate Representation Generator

Also you can share the yoga exercise to your friends via Email and SMS
RB PRP MD VB DT NN NN PRP NNS NNP NNP

share

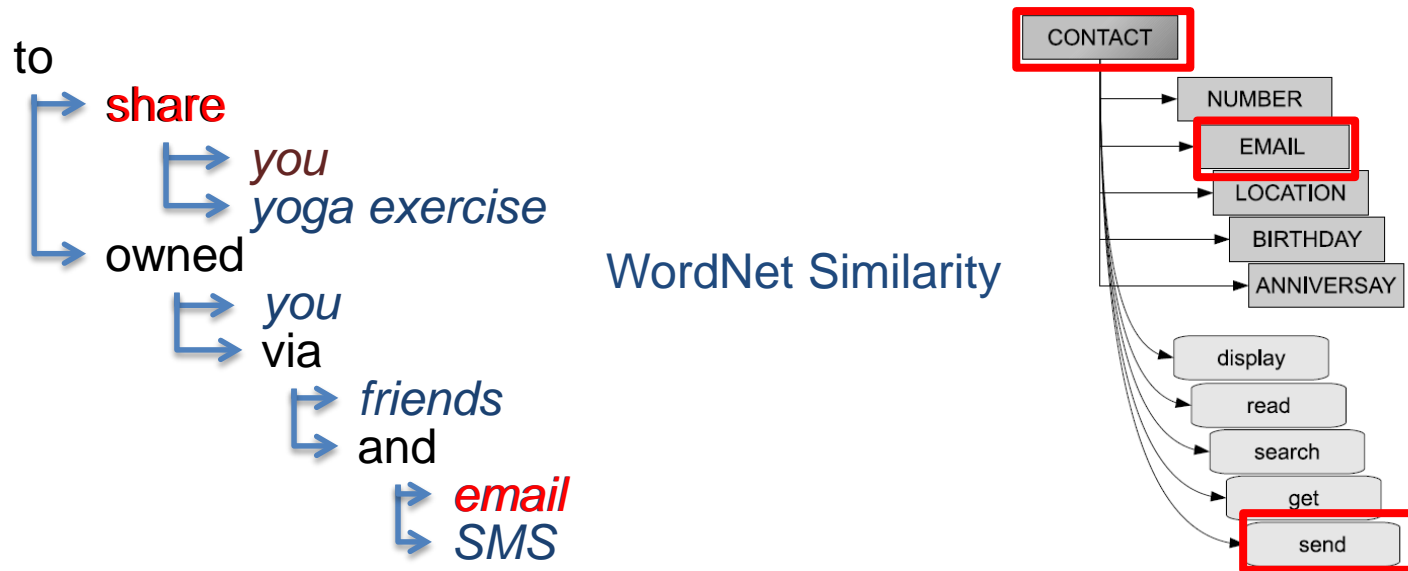


to

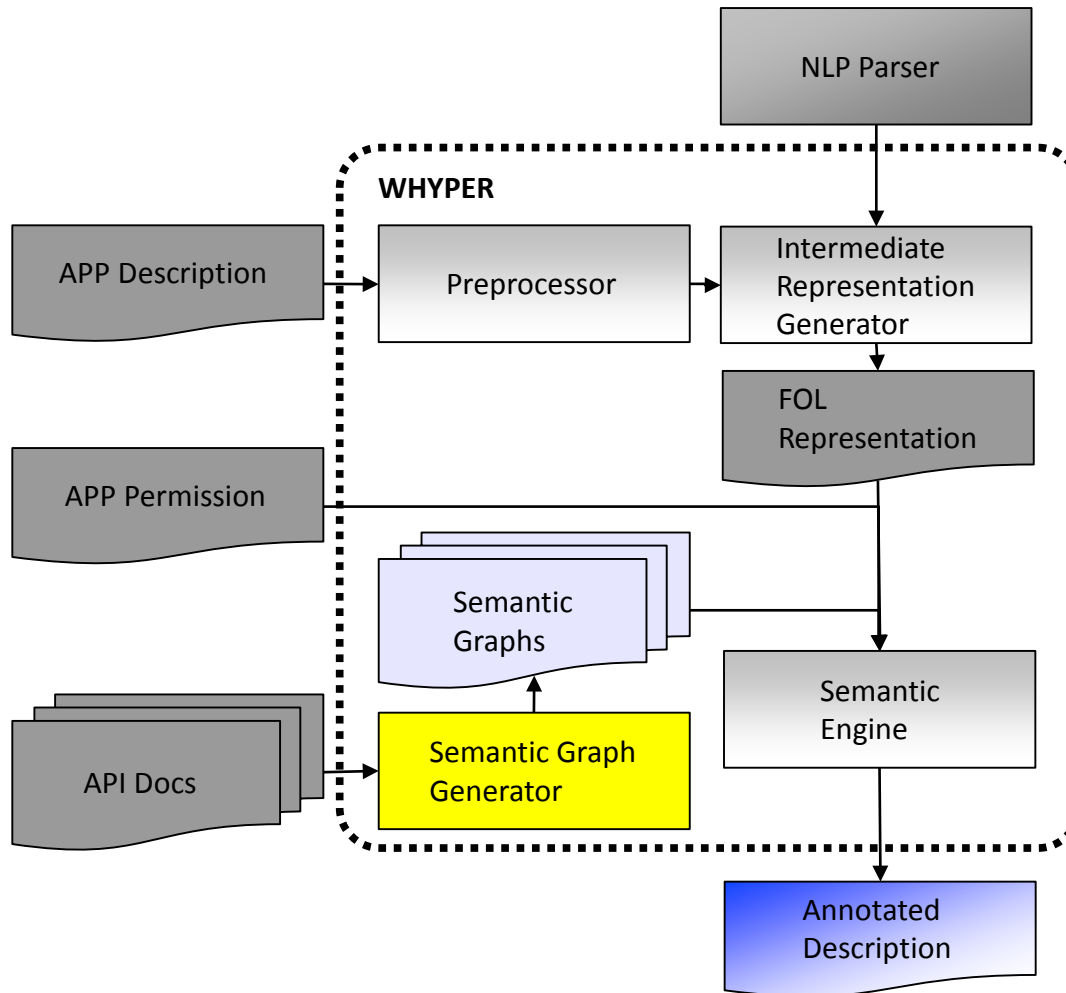


Semantic Engine

“Also you can share the yoga exercise to your friends via Email and SMS.”



WHYPER Framework



Semantic-Graph Generator

public static class



Summary: [Nested Classes](#) | [Constants](#) | [Inherited Constants](#) | [Fields](#) | [Methods](#) | [Inherited Methods](#) | [\[Expand All\]](#)

ContactsContract.Contacts

Added in API level 5

extends [Object](#)

implements [BaseColumns](#) [ContactsContract.ContactNameColumns](#) [ContactsContract.ContactOptionsColumns](#) [ContactsContract.ContactStatusColumns](#) [ContactsContract.ContactsColumns](#)

[java.lang.Object](#)

↳ [android.provider.ContactsContract.Contacts](#)

Class Overview

Constants for the contacts table, which contains a record per aggregate of raw contacts representing the same person.

Operations

Insert

A Contact cannot be created explicitly. When a raw contact is inserted, the provider will first try to find a Contact representing the same person. If one is found, the raw contact's `CONTACT_ID` column gets the `_ID` of the aggregate Contact. If no match is found, the provider automatically inserts a new Contact and puts its `_ID` into the `CONTACT_ID` column of the newly inserted raw contact.

Update

Only certain columns of Contact are modifiable: `TIMES_CONTACTED`, `LAST_TIME_CONTACTED`, `STARRED`, `CUSTOM_RINGTONE`, `SEND_TO_VOICEMAIL`. Changing any of these columns on the Contact also changes them on all constituent raw contacts.

Delete

Be careful with deleting Contacts! Deleting an aggregate contact deletes all constituent raw contacts. The corresponding sync adapters will notice the deletions of their respective raw contacts and remove them from their back end storage.

Query

- If you need to read an individual contact, consider using `CONTENT_LOOKUP_URI` instead of `CONTENT_URI`.
- If you need to look up a contact by the phone number, use `PhoneLookup.CONTENT_FILTER_URI`, which is optimized for this purpose.
- If you need to look up a contact by partial name, e.g. to produce filter-as-you-type suggestions, use the `CONTENT_FILTER_URI` URI.

Semantic-Graph Generator

- Systematic approach to infer graphs
 - Find related API documents using Pscout [CCS'12]
 - Identify resource associated with permissions from the API class name
 - `ContactsContract.Contacts`
 - Inspect the member variables and member methods to identify actions and subordinate resources
 - `ContactsContract.CommonDataKinds.Email`

Evaluation

- Subjects
 - Permissions: READ_CONTACTS, READ_CALENDAR, RECORD_AUDIO
 - 581/600* application descriptions (English only)
 - 9,953 sentences
- Research Questions
 - **RQ1**: What are the precision, recall, and F-Score of WHYPER in identifying permission sentences?
 - **RQ2**: How effective is WHYPER in identifying permission sentences, compared to keyword-based searching

Subject Statistics

Permissions	#N	#S	S_p
READ_CONTACTS	190	3,379	235
READ_CALENDAR	191	2,752	283
RECORD_AUDIO	200	3,822	245
TOTAL	581	9,953	763

Classification

- TP: WHYPER(s) && Manual(s)
- FP: WHYPER(s) && not(Manual(s))
- TN: not(WHYPER(s)) && not(Manual(s))
- FN: not(WHYPER(s)) && Manual(s)

RQ1 Results: Effectiveness

Permission	S_i	TP	FP	FN	TN	Prec.	Recall	F-Score	Acc
READ_CONTACT S	204	186	18	49	2,930	91.2	79.2	84.8	97.9
READ_CALENDAR	288	241	47	42	2,422	83.7	85.2	84.5	96.8
RECORD_AUDIO	259	195	64	50	3,470	75.3	79.6	77.4	97.0
TOTAL	751	622	129	141	9,061	82.8	81.5	82.2	97.3

- Out of 9,061 sentences, only 129 flagged as FPs
- Among 581 apps, 109 apps (18.8%) contain at least one FP
- Among 581 apps, 86 apps (14.8%) contain at least one FN

R2 Results: Comparison to Keyword-based search

Permission	Keywords
READ_CONTACTS	contact, data, number, name, email
READ_CALENDAR	calendar, event, date, month, day, year
RECORD_AUDIO	record, audio, voice, capture, microphone

Permission	Delta Precision	Delta Recall	Delta F-score	Delta Accuracy
READ_CONTACTS	50.4	1.3	31.2	7.3
READ_CALENDAR	39.3	1.5	26.4	9.2
RECORD_AUDIO	36.9	-6.6	24.3	6.8
WHYPER Improvement	41.6	-1.2	27.2	7.7

Results Analysis: False Positives

- Incorrect Parsing
 - “MyLink Advanced provides full synchronization of all Microsoft Outlook emails (inbox, sent, outbox and drafts), contacts, calendar, tasks and notes with all Android phones via USB”
- Synonym Analysis
 - “You can now **turn** recordings into ringtones.”

Results Analysis: False Negatives

- Incorrect parsing
 - Incorrect identification of sentence boundaries and limitations of underlying NLP infrastructure
- Limitations of Semantic Graphs
 - Manual Augmentation
 - Microphone (*blow into*) and call (*record*)
 - Significant improvement of delta recalls: **-6.6%** to **0.6%**
 - Future: automatic mining from user comments and forums

Broader Applicability

- Generalization to other permissions
 - User-understandable permissions: calls, SMS
 - Problem areas
 - Location and phone identifiers (widely abused)
 - Internet (nearly every app requires)

Dataset and Paper

- Our code and datasets are available at <https://sites.google.com/site/whypermission/>
- Rahul Pandita, Xusheng Xiao, Wei Yang, William Enck, and Tao Xie. **WHYPER: Towards Automating Risk Assessment of Mobile Applications.** *In Proc. 22nd USENIX Security Symposium (USENIX Security 2013)*
<http://www.enck.org/pubs/pandita-sec13.pdf>

Outline

- Introduction
- Background on text analytics
- Case Study 1: App Markets
- **Case Study 2: ACP Rules**
- Wrap-up

Access Control Policies (ACP)

- Access control is often governed by security policies called Access Control Policies (ACP)
 - Includes rules to control which principals have access to which resources

ex.

“The Health Care Personnel (HCP) does not have the ability to edit the patient's account.”

- A policy rule includes four elements
 - Subject – HCP
 - Action – edit
 - Resource - patient's account
 - Effect - deny

Access Control Vulnerabilities

eWEEK.COM

FaceBook Five Security

V

Facebook
permitted
without t

CWE and SANS Institute

Improper access control causes problems
(e.g., information exposures)

- Incorrect specification
- Incorrect enforcement

IT Security & Network Security News

3. Classic buffer overflow

4. Cross-site request forgery

5. **Improper access control (Authorization)**

6. ...

CWE Common Weakness Enumeration
A Community-Developed Dictionary of Software Weakness Types

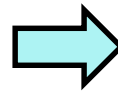
Problems of ACP Practice

- In practice, ACPs
 - Buried in requirement documents
 - Written in NL and not checkable
- NL documents could be large in size
 - Manual extraction is labor-intensive and tedious

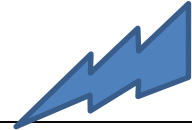
Overview of Text2Policy

Linguistic Analysis

A HCP should not change patient's account.



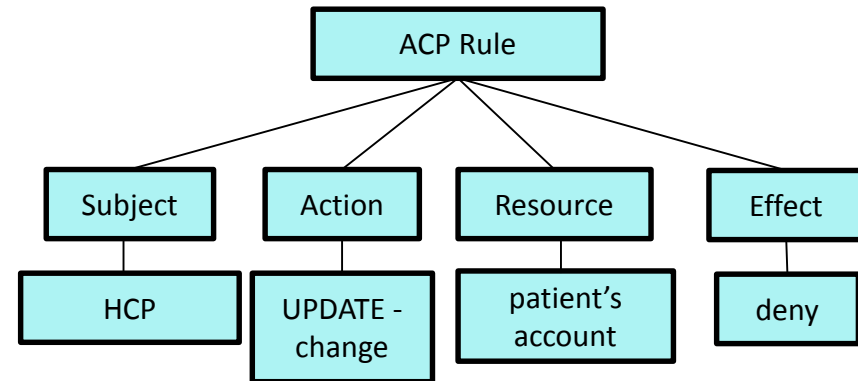
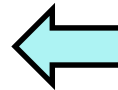
An [*subject*: HCP] should not [*action*: change] [*resource*: patient's account].



Model-Instance Construction



```
<Policy PolicyId="ACP2" RuleCombAlgId="deny-overrides">
  <Target/>
  <Rule Effect="Deny" RuleId="rule-1">
    <Target>
      <Subjects><Subject><SubjectMatch MatchId="string-equal">
        <AttrValue>HCP</AttrValue>
        <SubjectAttrDesignator.../></SubjectMatch></Subject>
      </Subjects>
      <Resources><Resource><ResourceMatch MatchId="string-equal">
        <AttrValue>patient.account</AttrValue>
        <ResourceAttrDesignator.../></ResourceMatch></Resource>
      </Resources>
      <Actions><Action><ActionMatch MatchId="string-equal">
        <AttrValue DataType="string">UPDATE</AttrValue>
        <ActionAttrDesignator.../></ActionMatch></Action>
      </Actions>
    </Target></Rule></Policy>
```



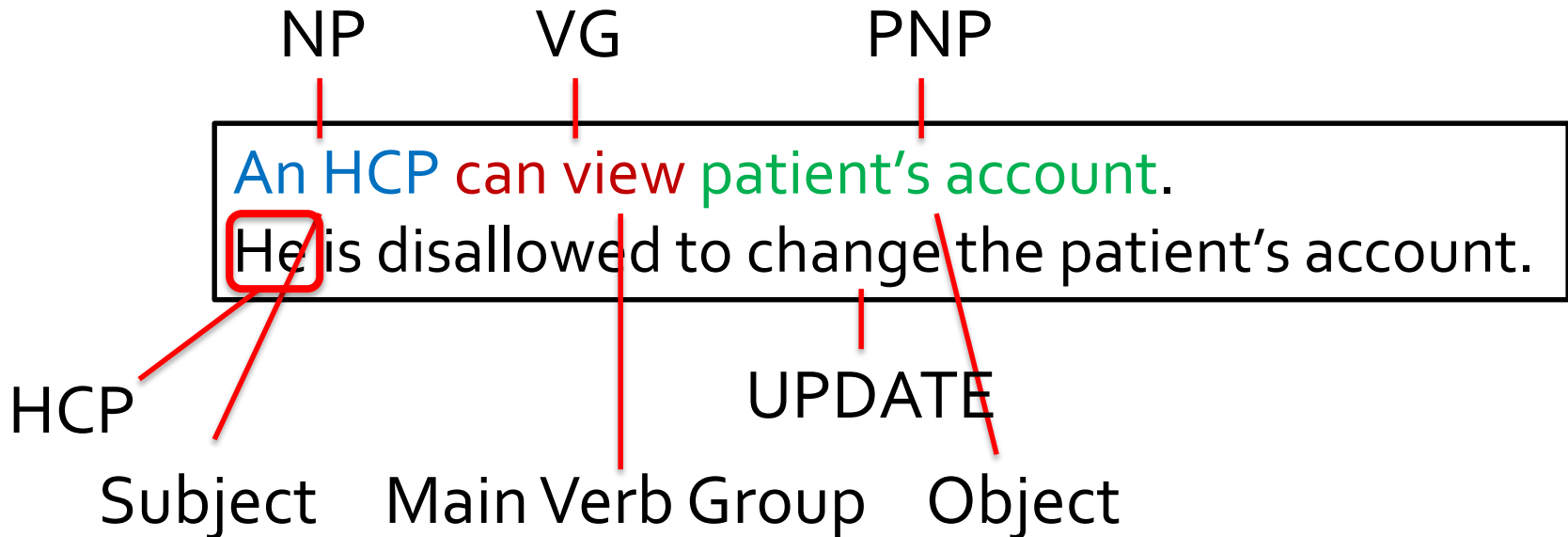
Transformation

Linguistic Analysis

- Incorporate syntactic and semantic analysis
 - **syntactic** structure -> noun group, verb group, etc.
 - **semantic** meaning -> subject, action, resource, negative meaning, etc.
- Provide New techniques for model extraction
 - Identify ACP and AS sentences
 - Infer semantic meaning

Common Techniques

- Shallow parsing
- Domain dictionary
- Anaphora resolution



Technical Challenges (TC) in ACP Extraction

ACP 1: An HCP cannot change patient's account.

ACP2: An HCP is disallowed to change patient's account.

- TC1: Semantic Structure Variance
 - different ways to specify the same rule
- TC2: Negative Meaning Implicitness
 - verb could have negative meaning

Semantic-Pattern Matching

- Address TC1 Semantic Structure Variance
- Compose pattern based on grammatical function

ex.

An HCP is disallowed to change the patient's account.

passive voice

followed by to-infinitive phrase

Negative-Expression Identification

- Address TC2 Negative Meaning Implicitness

- Negative expression

- “not” in subject:

- ex.

- No** HCP can edit patient’s account.

- “not” in verb group:

- ex.

- HCP can **not** edit patient’s account.

- HCP can **never** edit patient’s account.

- Negative meaning words in main verb group

- ex.

- An HCP is **disallowed** to change the patient’s account.**

AS: Syntactic-Pattern Matching

- Syntactic elements
 - Subject , Main verb, Object
- Subject and Object Checking
 - subject is a not a user or object is not a resource

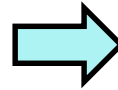
ex. **The prescription list** should include medication,
the name of the doctor. . .

- Filtering negative-meaning sentences
 - Negative sentences tend not to describe ASs

Overview of Text2Policy

Linguistic Analysis

A HCP should not change patient's account.

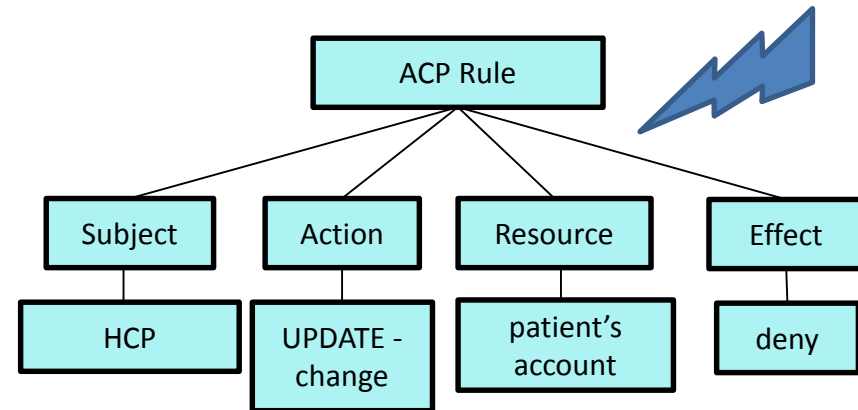
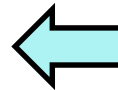


An [*subject*: HCP] should not [*action*: change] [*resource*: patient's account].

Model-Instance Construction



```
<Policy PolicyId="ACP2" RuleCombAlgId="deny-overrides">
  <Target/>
  <Rule Effect="Deny" RuleId="rule-1">
    <Target>
      <Subjects><Subject><SubjectMatch MatchId="string-equal">
        <AttrValue>HCP</AttrValue>
        <SubjectAttrDesignator.../></SubjectMatch></Subject>
      </Subjects>
      <Resources><Resource><ResourceMatch MatchId="string-equal">
        <AttrValue>patient.account</AttrValue>
        <ResourceAttrDesignator.../></ResourceMatch></Resource>
      </Resources>
      <Actions><Action><ActionMatch MatchId="string-equal">
        <AttrValue DataType="string">UPDATE</AttrValue>
        <ActionAttrDesignator.../></ActionMatch></Action>
      </Actions>
    </Target></Rule></Policy>
```



Transformation

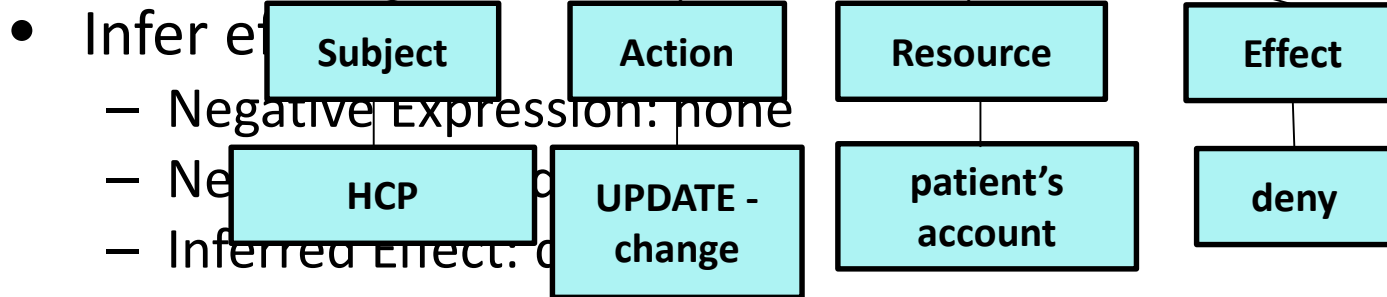
ACP Model-Instance Construction

ex. An HCP is disallowed to change the patient's account.

- Identify subject, action, and resource:

- Subject: HCP
- Action: change
- Resource: patient's account

ACP Rule



- Access Control Rule Extraction (ACRE) approach [ACSAC'14] discovers more patterns
 - Able to handle existing, unconstrained NL texts

Evaluation – RQs

- RQ1: How effectively does Text2Policy identify ACP sentences in NL documents?
- RQ2: How effectively does Text2Policy extract ACP rules from ACP sentences?

Evaluation – Subject

- iTrust open source project
 - <http://agile.csc.ncsu.edu/iTrust/wiki/>
 - 448 use-case sentences (37 use cases)
 - preprocessed use cases
- Collected ACP sentences
 - 100 ACP sentences
 - From 17 sources (published papers and websites)
- A module of an IBMApp (financial domain)
 - 25 use cases

RQ1 ACP Sentence Identification

- Apply Text2Policy to identify ACP sentences in iTrust use cases and IBMApp use cases

Subjects	# Sent.	# ACP Sent.	# Ident.	<i>FP</i>	<i>FN</i>	<i>Prec</i>	<i>Rec</i>	<i>F₁</i>
iTrust	448	117	119	16	14	86.6%	88.0%	87.3
IBMApp	479	24	23	0	1	100.0%	95.8%	97.9
Total	927	141	142	16	15	88.7%	89.4%	89.1

- Text2Policy effectively identifies ACP sentences with precision and recall more than 88%
- Precision on IBMApp use cases is better
 - proprietary use cases are often of higher quality compared to open-source use cases

Evaluation –

RQ2 Accuracy of Policy Extraction

- Apply Text2Policy to extract ACP rules from ACP sentences

Subjects	# ACP Sent.	# Extracted	Accu.
iTrust	217	187	86.2%
IBMAApp	24	21	87.5%
Total	241	208	86.3%

- Text2Policy effectively extracts ACP model instances with accuracy above 86%

Dataset and Paper

- Our datasets are available at <https://sites.google.com/site/asergroup/projects/text2policy>
- Xusheng Xiao, Amit Paradkar, Suresh Thummalapenta, and Tao Xie. **Automated Extraction of Security Policies from Natural-Language Software Documents.** In *Proc. 20th ACM SIGSOFT Symposium on the Foundations of Software Engineering (FSE 2012)*
<http://taoxie.cs.illinois.edu/publications/fse12-nlp.pdf>
- John Slankas, Xusheng Xiao, Laurie Williams, and Tao Xie. **Relation Extraction for Inferring Access Control Rules from Natural Language Artifacts.** In *Proc. 30th Annual Computer Security Applications Conference (ACSAC 2014)*
<http://taoxie.cs.illinois.edu/publications/acsac14-nlp.pdf>

Outline

- Introduction
- Background on text analytics
- Case Study 1: App Markets
- Case Study 2: ACP rules
- **Wrap-up**

Take-away

- Computing systems contain textual data that partially represents **expectation context**.
- Text analytics and natural language processing offers an opportunity to automatically extract that semantic context
 - Need to be careful in the security domain (e.g., social engineering)
 - But potential for improved security decisions

Future Directions

- Only beginning to study text analytics for security
 - Many sources of natural language text
 - Many unexplored domains
 - Use text analytics in software engineering as inspiration
 - <https://sites.google.com/site/text4se/>
- Hard problem: to what extent can we formalize “expectation context”?
- Creation of open datasets (annotation is time intensive)
- Apply to real-world problems



Thank you!



Questions?

Tao Xie

University of Illinois at Urbana-Champaign
<http://taoxie.cs.illinois.edu/>
taoxie@illinois.edu

William Enck

North Carolina State University
<http://www.enck.org>
enck@cs.ncsu.edu

Acknowledgment: We thank authors of the original slides that some slides from this tutorial were adapted from. The work is supported in part by NSA Science of Security Lablet grants, a Google Faculty Research Award, NSF grants CNS-1513939/1513690, CCF-1434596, CNS-1434582, CNS-1222680, and CNS-1253346.