

ภาคผนวก

ภาคผนวก ก

ตัวอย่างการเตรียมคลังข้อมูล

## New Antioxidant C-Glucosylxanthenes from the Stems of *Arrabidaea samydoides*

Patrícia Mendonça Pauletti,<sup>†</sup> Ian Castro-Gamboa,<sup>†</sup> Dulce Helena Siqueira Silva,<sup>†</sup> Maria Claudia Marx Young,<sup>‡</sup> Daniela Maria Tomazela,<sup>§</sup> Marcos Nogueira Eberlin,<sup>§</sup> and Vanderlan da Silva Bolzani<sup>\*,†</sup>

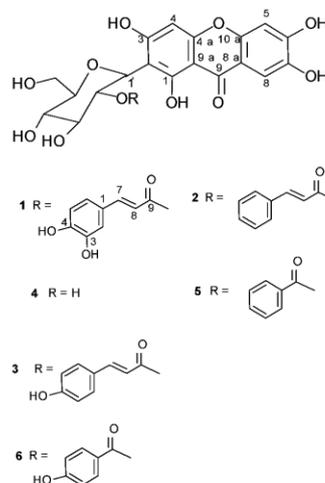
NuBBE- Núcleo de Biossíntese, Bioensaios e Ecofisiologia de Produtos Naturais, Instituto de Química, Universidade Estadual Paulista, UNESP, CP 355, 14801-970, Araraquara, SP, Brazil, Seção de Fisiologia e Bioquímica de Plantas, Instituto de Botânica, CP 4009, 01061-970, São Paulo, Brazil, and Thompson Mass Spectroscopy Laboratory, Instituto de Química, Universidade Estadual de Campinas, UNICAMP, CP 6154, 13083-970, Campinas, SP, Brazil

Received March 7, 2003

Three new C-glucosylxanthenes, 2-(2'-*O*-*trans*-caffeoyl)-C- $\beta$ -D-glucopyranosyl-1,3,6,7-tetrahydroxyxanthone (1), 2-(2'-*O*-*trans*-cinnamoyl)-C- $\beta$ -D-glucopyranosyl-1,3,6,7-tetrahydroxyxanthone (2), and 2-(2'-*O*-*trans*-coumaroyl)-C- $\beta$ -D-glucopyranosyl-1,3,6,7-tetrahydroxyxanthone (3), were isolated from the stems of *Arrabidaea samydoides*, in addition to three known C-glucosylxanthenes, mangiferin (4), 2-(2'-*O*-benzoyl)-C- $\beta$ -D-glucopyranosyl-1,3,6,7-tetrahydroxyxanthone (5), and muraxanthone (6). Their chemical structures were assigned on the basis of MS and 1D and 2D NMR experiments. Xanthenes 1–6 showed moderate free radical scavenging activity against 1,1-diphenyl-2-picrylhydrazyl (DPPH) as well as antioxidant activity evidenced by redox properties measured on EICD-HPLC.

As part of our bioprospecting program Biota-FAPESP (The Virtual Institute of Biodiversity), whose main goal is to discover potential antitumoral, antifungal, and antioxidant agents from plants of the Cerrado and Atlantic forest, we have screened hundreds of plants collected in the State of São Paulo. Among these, *Arrabidaea samydoides* was chosen for detailed chemical investigation due to prior antioxidant activity revealed on a TLC autographic assay sprayed with  $\beta$ -carotene solution, and to our knowledge there are no previous reports on chemical and biological studies. This species belongs to the family Bignoniaceae, which contains about 120 genera and 800 species distributed throughout tropical regions of South America and Africa.<sup>1</sup> Species from the genus *Arrabidaea* have been used in traditional medicine for wound asepsis and treating intestinal disorders.<sup>2</sup> In northeast Brazil, *Arrabidaea chica* is used in tattoos by Indians due to the pigments carajurin and carajurone.<sup>2,3</sup> A literature review indicated that this genus is a source of anthocyanins, flavonoids, and tannins.<sup>3–7</sup> The ethanolic extract from the stems showed promising antioxidant activity and led to the isolation of three new C-glucopyranosylxanthenes, 2-(2'-*O*-*trans*-caffeoyl)-C- $\beta$ -D-glucopyranosyl-1,3,6,7-tetrahydroxyxanthone (1), 2-(2'-*O*-*trans*-cinnamoyl)-C- $\beta$ -D-glucopyranosyl-1,3,6,7-tetrahydroxyxanthone (2), and 2-(2'-*O*-*trans*-coumaroyl)-C- $\beta$ -D-glucopyranosyl-1,3,6,7-tetrahydroxyxanthone (3), and the known mangiferin (4),<sup>8</sup> 2-(2'-*O*-benzoyl)-C- $\beta$ -D-glucopyranosyl-1,3,6,7-tetrahydroxyxanthone (5),<sup>9</sup> and muraxanthone (6).<sup>10</sup> In this paper, we report the isolation, structure elucidation, and antioxidant properties of these C-glucopyranosylxanthenes.

Compounds 4 and 6 were identified by comparison with previously published NMR and other physical data.<sup>8,10</sup> Compound 5 was described previously as a mixture of three isomers from *Hymenophyllum recurvum*.<sup>9</sup> Only a few <sup>13</sup>C NMR data were analyzed and discussed. In this paper we



describe the complete <sup>1</sup>H, <sup>13</sup>C NMR and ES-MS/MS data for this compound.

Compound 1 was shown to have the molecular formula C<sub>28</sub>H<sub>28</sub>O<sub>14</sub> [M - H]<sup>-</sup> m/z 583.1008, by analysis of the negative-ion HRESIMS. The IR spectrum showed bands at 3370, 1615, and 1474 cm<sup>-1</sup> accounting for hydroxyl, conjugated carbonyl, and aromatic groups, respectively. The <sup>13</sup>C NMR spectrum showed six signals for hydroxymethine carbons, suggesting the presence of a sugar moiety, and 22 signals for sp<sup>2</sup> carbons, which could be assigned to three aromatic rings, and also two carbonyls and one additional olefinic function. In the <sup>1</sup>H NMR spectrum (Table 1) of 1, a caffeoyl moiety was identified by signals at  $\delta$  6.79 (1H, d, *J* = 2.0 Hz, H-2'), 6.58 (1H, d, *J* = 8.0 Hz, H-5'), and 6.67 (1H, br d, *J* = 8.0 Hz, H-6'),

\* Author to whom correspondence should be addressed. Tel: 55(16)-2016660. Fax: 55(16)2227932. E-mail: bolzani@iq.unesp.br.

<sup>†</sup> Instituto de Química, Universidade Estadual Paulista-UNESP.  
<sup>‡</sup> Seção de Fisiologia e Bioquímica de Plantas, Instituto de Botânica.  
<sup>§</sup> Instituto de Química, Universidade Estadual de Campinas-UNICAMP.

ภาพตัวอย่างเอกสารที่นำมาใช้เตรียมคลังข้อมูล

## ตัวอย่างข้อมูลที่ทำการ annotation และนำเอา stop word ออก

<stuff\_relation class="yes">Three new C-glucosylxanthenes, 2-(2'-O-trans-caffeoyl)-C-β-D-glucopyranosyl-1,3,6,7-tetrahydroxyxanthone(**1**), 2-(2'-O-trans-cinnamoyl)-C-β-D-glucopyranosyl-1,3,6,7-tetrahydroxyxanthone 2-(2'-O-trans-coumaroyl)-C-β-D-glucopyranosyl-1,3,6,7-tetrahydroxyxanthone isolated stems *Arrabidaea samydoidea* C-glucosylxanthenes, mangiferin 2-(2'-O-benzoyl)-C-β-D-glucopyranosyl-1,3,6,7-tetrahydroxyxanthone muraxanthone</stuff\_relation>

<stuff\_relation class="no">chemical structures assigned MS 1D 2D NMR experiments</stuff\_relation>

<stuff\_relation class="no">Xanthenes **1-6** showed moderate radical scavenging 1,1-diphenyl-2-picrylhydrazyl (DPPH) antioxidant evidenced redox properties measured EICD-HPLC</stuff\_relation>

<stuff\_relation class="no">bioprospecting program Biota-FAPESP Virtual Institute Biodiversity goal discover potential antitumoral antifungal antioxidant agents Cerrado Atlantic forest screened hundreds collected State São Paulo</stuff\_relation>

<stuff\_relation class="no">*Arrabidaea samydoidea* chosen detailed chemical investigation prior antioxidant revealed TLC autographic assay sprayed α-carotene solution knowledge reports chemical biological studies</stuff\_relation>

<stuff\_relation class="no">species belongs family Bignoniaceae contains 120 genera 800 species distributed tropical regions South America Africa</stuff\_relation>

<stuff\_relation class="no">Species genus *Arrabidaea* traditional medicine wound asepsis treating intestinal disorders</stuff\_relation>

<stuff\_relation class="no">northeast Brazil *Arrabidaea chica* tattoos Indians pigments carajurin carajurone</stuff\_relation>

<stuff\_relation class="no">literature review indicated genus source anthocyanins flavonoids tannins.</stuff\_relation>

<stuff\_relation class="no">ethanolic extract stems showed promising antioxidant isolation C-glucopyranosylxanthenes 2-(2'-O-trans-caffeoyl)-C-α-D-glucopyranosyl-1,3,6,7-tetrahydroxyxanthone 2-(2'-O-trans-cinnamoyl)-C-α-D-glucopyranosyl-1,3,6,7-tetrahydroxyxanthone 2-(2'-O-trans-coumaroyl)-C-α-D-glucopyranosyl-1,3,6,7-tetrahydroxyxanthone mangiferin 2-(2'-O-benzoyl)-C-α-D-glucopyranosyl-1,3,6,7-tetrahydroxyxanthone muraxanthone</stuff\_relation>

<stuff\_relation class="no">paper report isolation structure elucidation antioxidant properties C-glucopyranosylxanthenes</stuff\_relation>

<stuff\_relation class="no">Compounds **4 6** identified comparison published NMR physical data</stuff\_relation>

<stuff\_relation class="no">Compound **5** described mixture isomers *Hymenophyllum recurvum*</stuff\_relation>

<stuff\_relation class="no"><sup>13</sup>C NMR data analyzed discussed</stuff\_relation>

<stuff\_relation class="no">paper describe complete <sup>1</sup>H, <sup>13</sup>C NMR ES-MS/MS data compound</stuff\_relation>

<stuff\_relation class="no">Compound **1** shown molecular formula C<sub>28</sub>H<sub>23</sub>O<sub>14</sub> [M - H]<sup>-</sup>. m/z 583.1008 analysis negative-ion HRESIMS</stuff\_relation>

<stuff\_relation class="no">IR spectrum showed bands 3370 1615 1474 cm<sup>-1</sup> accounting hydroxyl conjugated carbonyl aromatic groups</stuff\_relation>

<stuff\_relation class="no"><sup>13</sup>C NMR spectrum showed signals hydroxymethine carbons suggesting presence sugar moiety 22 signals sp<sup>2</sup> carbons assigned aromatic rings carbonyls additional olefinic function</stuff\_relation>

ภาคผนวก ข

ตัวอย่างประโยคที่เป็นความสัมพันธ์แบบสต๊าฟฟ์ที่สกัดได้

id	ตัวอย่างประโยชน์ที่เป็นความสัมพันธ์แบบสตัดฟ์ที่สกัดได้
88	<i>C</i> -glucosylxanthenes, 2-(2'- <i>O</i> -trans-caffeoyl)- <i>C</i> - <i>O</i> - <i>D</i> -glucopyranosyl-1,3,6,7-tetrahydroxyxanthone 2-(2'- <i>O</i> -trans-cinnamoyl)- <i>C</i> - <i>O</i> - <i>D</i> -glucopyranosyl-1,3,6,7-tetrahydroxyxanthone 2-(2'- <i>O</i> -trans-coumaroyl)- <i>C</i> - <i>O</i> - <i>D</i> -glucopyranosyl-1,3,6,7-tetrahydroxyxanthone isolated stems <i>Arrabidaea samydoides</i>
217	5- <i>epi</i> -vibsanin G 18- <i>O</i> -methylvibsanin G vibsanin M aldovibsanin C isolated acetone extract leaves flowers <i>Viburnum odoratissimum</i>
318	(8 <i>R</i> *)-8-bromo-10- <i>epi</i> - $\beta$ -snyderol (8 <i>S</i> *)-8-bromo- $\beta$ -snyderol 5-bromo-3-(3'-hydroxy-3'-methylpent-4'-enylidene)-2,4,4-trimethylcyclohexanone isolated chloroform-methanol extract <i>Laurencia obtuse</i>
811	jujubogenin 3- <i>O</i> - <i>R</i> - <i>L</i> -arabinofuranosyl(1 <i>f</i> 2)-[ $\beta$ - <i>D</i> -glucopyranosyl(1 <i>f</i> 6) <i>O</i> - <i>D</i> -glucopyranosyl(1 <i>f</i> 3)]- <i>R</i> - <i>L</i> -arabinopyranoside jujubogenin 3- <i>O</i> - <i>R</i> - <i>L</i> -arabinofuranosyl(1 <i>f</i> 2)-{6- <i>O</i> -[3-hydroxy-3-methylglutaryl]- <i>O</i> - <i>D</i> -glucopyranosyl(1 <i>f</i> 3)}- <i>R</i> - <i>L</i> -arabinopyranoside <i>O</i> -hydroxylup-20(29)-en-27,28-dioic acid 28- <i>O</i> - $\beta$ - <i>D</i> -glucopyranosyl(1 <i>f</i> 2)-[ <i>O</i> - <i>D</i> -xylopyranosyl(1 <i>f</i> 3)]- <i>O</i> - <i>D</i> -xylopyranosyl(1 <i>f</i> 2)- $\beta$ - <i>D</i> -glucopyranoside ester jujubogenin 3- <i>O</i> - <i>R</i> - <i>L</i> -arabinofuranosyl(1 <i>f</i> 2)-[ $\beta$ - <i>D</i> -glucopyranosyl(1 <i>f</i> 3)]- <i>R</i> - <i>L</i> -arabinopyranoside <i>O</i> -hydroxylup-20(29)-ene-27,28-dioic acid isolated methanol extract stems <i>Anomospermum grandifolium</i> .
878	5 <i>R</i> ,9 <i>R</i> ,10 $\beta$ ,13 <i>R</i> -tetraacetoxy-14 $\beta$ - <i>O</i> -( $\beta$ - <i>D</i> -glucopyranosyl)taxa-4(20),11-diene 1 $\beta$ ,2 <i>R</i> ,9 <i>R</i> ,10 $\beta$ -tetrahydroxy-5 <i>R</i> -cinnamoyoxytaxa-4(20),11-dien-13-one (2), 2 <i>R</i> ,9 <i>R</i> ,10 $\beta$ -trihydroxy-5 <i>R</i> -cinnamoyoxytaxa-4(20),11-dien-13-one 9 <i>R</i> -acetoxy-2 <i>R</i> ,10 $\beta$ -dihydroxy-5 <i>R</i> -cinnamoyoxytaxa-4(20),11-dien-13-one 2 <i>R</i> ,10 $\beta$ -diacetoxy-1 $\beta$ ,9 <i>R</i> -dihydroxy-5 <i>R</i> -cinnamoyoxy-3,11-cyclotaxa-4(20)-dien-13-one identified <i>Taxus baccata</i>
1188	wilfordine alata mine wilforidine alatusinine euonine euonymine ebenifoline forrestine mayteine 4-hydroxy-7- <i>epi</i> -chuchuhuanine isolated leaves <i>Maytenus chiapensis</i> .
1364	<i>Asparagus cochichinensis</i> led isolation asparacoside asparacosins A 3''-methoxyasparenediol 3'-hydroxy-4'-methoxy-4'-dehydroxy nyasol asparenediol nyasol 3''-methoxy nyasol 1,3-bis-di- <i>p</i> -hydroxyphenyl-4-penten-1-one trans-coniferyl alcohol

1734	C-glycosides 2''-O-(2'''-methylbutyryl)isowertisin 3''-O-(2'''-methylbutyryl)-isowertisin 2''-O-(2'''-methylbutyryl)vitexin 2''-O-(2'''-methylbutyryl)orientin 2''-O-(3''',4'''-dimethoxybenzoyl)vitexin 2''-O-(3''',4'''-dimethoxybenzoyl)orientin isowertisin isolated flowers <i>Troliius ledebouri</i> .
2163	2,3,6,8-tetrahydroxy-1-(3-methylbut-2-enyl)-5-(2-methylbut-3-en-2-yl)-9H-xanthen-9-one isolated root bark <i>Cudrania</i>
2205	5R,7R,10 $\beta$ H-3-patchoulen-2-one 5R,7R10 $\beta$ H-4(14)-patchoulen-2R-ol 9R,10 $\beta$ -dihydroxy-2 $\beta$ ,4 $\beta$ -peroxy-1R,5 $\beta$ ,7RH-guaiane isolated aerial parts <i>Croton arboreous</i>
2762	10,11-dimethoxynareline alstohentine alstomicine 16-hydroxyalstonisine 16-hydroxyalstonal 16-hydroxy-N(4)-demethylalstophyllal oxindole alstophyllal 6-oxoalstophylline 6-oxoalstophyllal ones obtained leaf extract Malayan <i>Alstonia macrophylla</i> .
2772	14 $\beta$ -benzoyloxybaccatin IV 14 $\beta$ -benzoyloxy-13-deacetyl baccatin IV 14 $\beta$ -benzoyloxy-2-deacetyl baccatin VI isolated leaves stems <i>Taxus chinensis</i> .
3846	4 $\beta$ -hydroxy-19-normanoyl oxide 4R-hydroxy-18-normanoyl oxide 18-O-R-L-arabinopyranosylmanoyl oxide jhanol 18-hydroxy-13-epi-manoyl oxide isolated constituents Argentine collection <i>Grindelia scorzonifolia</i>

## **ภาคผนวก ก**

บทความการประชุม 7th International Conference on Computer Sciences and Convergence  
Information Technology (ICCIT2012) ณ กรุงโซล สาธารณรัฐเกาหลีใต้

# Automatic Stuff Relation Extraction from Scientific Documents for Natural Product Ontology Construction

Suriyasak Lertsakunsomboon  
 Search Engines and Intelligent Information Systems  
 Research Laboratory  
 Graduate Program in Web Engineering  
 Faculty of Information Technology  
 Dhurakij Pundit University  
 Bangkok, Thailand  
 535159090014@mydpu.net

Chaveevan Pechsiri  
 Search Engines and Intelligent Information Systems  
 Research Laboratory  
 Graduate Program in Web Engineering  
 Faculty of Information Technology  
 Dhurakij Pundit University  
 Bangkok, Thailand  
 chaveevan.pec@dpu.ac.th

**Abstract**— To extract Part-Whole relations, especially the stuff relation, from unstructured textual data is the challenging work. This paper presents how to automatically extract the stuff relation from technical documents on the Web for supporting chemical industries. The research extracts the stuff relation without applying POS (Part-of-Speech) annotation. There are three problems of extracting the stuff relation: a) the identification of stuff relation without POS annotation problem, b) the chemical-formula-embedded name entity determination problem and c) the genus-species name entity determination problem. We propose using Naive Bayes to learn the stuff relation. The results from our proposed methodology are 87% precision and 61% recall.

**Keywords**—stuff relation; scientific name entity; chemical name entity

## I. INTRODUCTION

Through out history, there are consistent amount of interests to extract chemicals of natural products for using in different areas of the pharmaceutical industry [1], the alternative energy research [2], and the commercial biorefinery system [3]. Therefore, several natural product researches work on analysis of natural substances of land and sea and of plants, microbes and animals where the research results bring significant benefits to the industries, especially the pharmaceutical industry. For example, many natural product researches focus on identifying novel bioactive constituents from mushrooms and medicinal plant. Recent studies identified antibacterial and cytotoxic mushroom species (where cytotoxic mushroom uses as anti-cancer medicines), as well as novel compounds from Bangladeshi medicinal plants. Since numerous scientific literatures exist on natural substances, the allocation of large amounts of time and resources are required by industries to seek source organism alternatives for a specific constituent. Therefore this research proposes to automatically extract the Part-Whole relation (i.e. "X is the component of Y") from the research documents to reveal the source organism substances of the natural product. According to M.E. Winston et. al. 1987, the

Part-Whole relation can be classified into six types: Component-Integral object (wheel-car), Member-Collection (soldier-army), Portion-Mass (meter-kilometer), Stuff-Object (alcohol-wine), Feature-Activity (paying-shopping), and Place-Area (oasis-desert). Therefore, this research aims to determine and extract the Part-Whole relation, especially Stuff-Object type, from the documents. Unlike the Component-Integral relation, the Stuff-Object relation (or Stuff relation) refers to a relation where the stuff cannot be physically separated from the object without altering its identity [4]. The Stuff relation is required for automatically constructing the natural product Ontology used to represent all natural product knowledge.

There are several researches work on the relation extraction [5][6][7][8], but literature on the Part-Whole (or Part-of) relation extraction is still lacking. Most of these researches worked on the text file involved with the linguistic patterns or the linguistic rule bases at the phrase level (e.g. "the oil and vinegar salad dressing") or the sentences level (e.g. "A vinaigrette-type salad dressing in United States and Canadian cuisine consists of water, vinegar or lemon juice, vegetable oil, chopped bell peppers corn syrup, and a blend of various herbs and spices."). Our research concerns only the sentence level since the Stuff relation expression on our corpus of the natural product research papers mostly occurs at the sentence level. In addition, our corpus contains several characteristics differ from other general corpora, especially the name entity (especially IUPAC nomenclature or chemical name entity) such as the chemical-formula-embedded name entity as shown the following with the underline.

*"Seven new lanostane-type triterpenes, hypo-crellois A-G (1-7), and six new hopane-type triterpenes, 7 $\beta$ ,15 $\alpha$ -dihydroxy-22(29)-hopene (8), 3 $\beta$ ,7 $\beta$ -dihydroxy-22(29)-hopene (9), 3 $\beta$ -acetoxyl-15 $\alpha$ -hydroxy-22(29)-hopene (10), 3 $\beta$ ,7 $\beta$ ,15 $\alpha$ ,22-tetrahydroxyhopane (11),*

*3 $\beta$ -acetoxy-7 $\beta$ ,15 $\alpha$ ,22-trihydroxyhopane (12), and 7 $\beta$ ,15 $\alpha$ ,22-trihydroxy-hopane (13), were isolated from the scale insect pathogenic fungus *Hypocrella* sp. BCC 14524."*

Where each number in the parenthesis except (29) stands for the identification number of the element in front of the parenthesis, for example:

"1-7" (1 to 7) stand for "*lanostane-type triterpenes, hypocrellols A-G*"

"8" stands for "*hopane-type triterpenes, 7 $\beta$ ,15 $\alpha$ -dihydroxy-22(29)-hopene*"

Another problem is the chemical name entity contains several commas whose functions differ from the general comma function separating word in a series (www.brighthubeducation.com). For example:

"....., 3,4-dihydroxy-9-methoxypterocarpan (vesticarpan) (5), 2',4,4'-trihydroxychalcone (isoliquiritigenin), and 7,4'-dihydroxyflavanone (liquiritigenin) (6), were isolated from the heartwood of *Platymiscium floribundum*."

where "2',4,4'-trihydroxychalcone" is one word. Whereas "dogs" / "cats" / "mice" is a one word in the following sentence.

"Dogs, cats, and mice are mammals."

All of these characteristics are involved in the three main problems; first is how to identify the Stuff relation. Second is how to determine the chemical formula name entity. Third is how to determine the organism (species) name entity. From all of these problems, we propose applying the Naïve Bayes machine learning technique to learn the stuff relation from a sentence that contains the interesting word set {"obtain", "isolate", "extract", ...}(from corpus behavior study) along with the stuff word set {"7,4'-dihydroxy-3'-methoxyisoflavone", "isoliquiritigenin", "(R)-4'-methoxydalbergione", ...}(from NBIC-PubChem database) and the object word set {"*Dalbergia louvelii*", "*Dendrolobium lanceolatum*", ...}(from NBIC-Taxonomy database).

In section II, related work is summarized. Problems in extracting the Stuff relation from the published research papers is described in section III and in section IV is purposed our framework for the Stuff relation extraction. In section V, we evaluate and conclude our proposed model.

## II. RELATED WORKS

There are some previous researches including [5][6][7][8] working on the relation extracting from texts as described in the following.

R. Girju et. al.(2003)'s work [5] is to present a ID3 (C4.5) learning technique for learning semantic constraints based on the lexico syntactic pattern, NP1 verb NP2 where NP1 and NP2 are noun phrases, to detect Part-Whole relations (meronymy) from the LA Times articles of TREC 9 text

collections. They also stated that the Part-Whole relations in WordNet were classified into three basic types: Member-of (e.g., UK IS-MEMBER-OF NATO), Stuff-of (e.g., carbon IS-STUFF-OF coal), and Part-of (e.g., leg IS-PART-OF table). According to their proposed model based on sentence level, the Part Whole relations are detected with an accuracy of 83% precision 98% recall from 10000 sentences.

In 2006, P. Pantel and M. Pennacchiotti's work [6] is to present the weakly-supervised algorithm, named Espresso, using the generic patterns to extract the semantic relations, especially the Is-a relation and the Part-of relation. The Espresso was applied to the very large corpora downloaded from webs with the generic pattern which was determined the pattern and instance reliability. The reliable patterns resulted in having high precision but often very low recall (e.g., "X consists of Y" for the Part-of relation). Their experimental results showed that their generic patterns substantially increased system recall with small effect on overall precision. The results of their system performance using the CHEM corpus with the Part-of relation is 51% precision and 46 relative recall whereas using TREC-9 with the Part-of relation is 70% precision 577 relative recall.

In 2007, K.Fundel et. al. [7] developed RelEx, an approach for relation extraction from free text, especially biomedical publications from million MEDLINE abstracts. RelEx was based on natural language preprocessing producing dependency parse trees dealing with gene and protein relations. The RelEx model involved the detection of co-occurrences of entities within sentences or abstracts and uses a small set of simple rules, applied for part-of-speech-tagging, noun-phrase-chunking and dependency. Their Relation Extraction is based on three linguistic rules frequently used in English language for describing relations where one rule is based on the sentence level and the other two rules are based on the phrase level as follow:

- (1) effector-relation-effectee (' $\alpha$  activates  $\beta$ ')
- (2) relation-of-effectee-by-effector ('Activation of  $\alpha$  by  $\beta$ ')
- (3) relation-between-effector-and-effectee('Interaction between  $\alpha$  and  $\beta$ ').

Finally, RelEx is estimated performance of both 80% precision and 80% recall.

G. I. Brown (2011) [8] presented an analysis of a relation extraction system using a support vector machine (SVM) classifier to the J.D. Power and Associates Sentiment Corpus separated into three style documents: professionally written reviews, blog reviews, and social networking reviews. The SVM features includes the word features involved the head noun of the phrase existing the least deep in the dependency parse tree, the entity types, and the token class. However, his research aims to study how the extraction system works on different styles of documents. The results of the Part-of relation extraction are 46%precision on average and 33% recall on average. Then, [8] concluded that the relation extraction task was being negatively impacted by the relation

classification itself and the poor tokenization or parsing of the documents.

However, all of the relation extraction techniques from [5][6][7][8] cannot be applied to our research. The complicated chemical-formula name entity and the organism name entity limit the use of POS standard tools (<http://open.xerox.com/Services/fst-nlp-tools/Consume/181>). Therefore, we apply Naïve Bayes to learn the Stuff relation from a sentence that contains the interesting-stem word set along with the Stuff word set and the object word set.

### III. PROBLEM OF STUFF RELATION EXTRACTION

There are three main problems: how to identify the Stuff relation, how to determine the chemical-formula-embedded name boundary, and how to determine the scientific name boundary.

#### A. How to identify Stuff Relation without POS annotation

In order to identify the Stuff relation without the POS annotation, we propose using the A-B-C or C-B-A sequence (where A is the Stuff word set, B is the interesting word set, and C is the object word set) occurring within one sentence to identify the Stuff relation using Naïve Bayes. A and C are obtained from NBIC-PubChem and NBIC-Taxonomydatabase (<http://pubchem.ncbi.nlm.nih.gov>) respectively (see Fig. 1).

“... Four new flavonoids (1-4), along with 13 known compounds, were isolated from the heartwood of *Dalbergia louvelii*. ...”

Fig. 1. Example of the linguistic expression in scientific documents.

where A is “flavonoids (1-4)”, B is “isolate”, and C is “*Dalbergia louvelii*”.

#### B. How to determine the chemical-formula-embedded name boundary

The chemical formula is very complex word (as shown in TABLE I.) where the chemical names are embedded within another chemical name (see section I). Therefore we used NBIC-PubChem to determine the chemical name.

TABLE I. LIST OF CHEMICAL FORMULA EXAMPLE.

Chemical Formula
<b>Example 1 :</b> (R)-4"-methoxydalbergione
<b>Example 2 :</b> 3-(2,4-dihydroxy-5-methoxy)phenyl-7-hydroxycoumarin
<b>Example 3 :</b> (R)-4"-methoxydalbergione
<b>Example 4 :</b> (7S,8R,1'S,5'S,6'R)- $\Delta^{2,8}$ -5',6'-dihydroxy-3'-methoxy-3,4-methylenedioxy-4'-oxo-8,1',7,5'-neolignan
<b>Example 5 :</b> 2,4-dimethoxy-5,6-methylenedioxy-1-(2-propenyl)benzene
<b>Example 6 :</b> 2'-hydroxy-6,4',6'',4'''-tetramethoxy-[7-O-7'']-bisiso flavone

Also, the chemical formula often contains “,” as separating word (see section I). We propose using NBIC-PubChem database to determine the chemical formula name entity.

#### C. How to determine the scientific (species and genus) name boundary

Often the species and genus names of species are written in multiple word italic form (as shown in TABLE II.). However, when the surrounding texts are also in italics, the identification of the scientific species names become a challenging task.

TABLE II. LIST OF SCIENTIFIC NAME EXAMPLE.

scientific name
<b>Example 1 :</b> <i>Dalbergia louvelii</i>
<b>Example 2 :</b> <i>Dendrolobium lanceolatum</i>
<b>Example 3 :</b> <i>Hypericum perforatum</i>
<b>Example 4 :</b> <i>Platymiscium floribundum</i>
<b>Example 5 :</b> <i>Dalbergia candanensis</i>

Therefore we propose using NBIC-Taxonomy database to solve this problem.

### IV. METHODOLOGY

There are 4 major steps including corpus preparation, name entity determination, Stuff relation learning, and Stuff relation extraction as shown in Fig. 2.

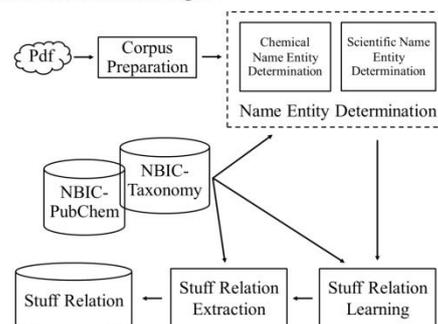


Fig. 2. System architecture.

#### A. Corpus Preparation

In the corpus preparation step, English scientific abstracts and introductions in chemical related areas are downloaded from online journals (20,000 sentences from 500 scientific documents in pdf format). The pdf documents are converted to text documents using PDFTextStream (<http://snowtide.com>). The corpus is separated into 2 parts; one part is 15,000 sentences for learning the stuff relation and the other part of 5,000 sentences for stuff relation extraction. The sentences with stuff relation are manually annotated as “yes” class and the others are “no” class (see Fig. 3.). Then the stop words are filtered out while all symbols, e.g. “,” “-” ... etc., still exist in the corpus.

```

<stuff_relation class="yes">Four new flavonoids (1-4), along with
13 known compounds, were isolated from the heartwood of
Dalbergia louvelii by following their potential to inhibit in vitro
the growth of Plasmodium falciparum.</stuff_relation>
<stuff_relation class="no">Of these, the ethyl acetate extract
obtained from the heartwood of Dalbergia louvelii R. Viguier
(Fabaceae).</stuff_relation>
<stuff_relation class="yes">Although several isoflavonoids have
been obtained from roots of P. floribundum, none of the
abovementioned compounds have been isolated previously from
this species.</stuff_relation>
<stuff_relation class="no">The cytotoxicity of the isolates
obtained herein from P. floribundum has been evaluated against a
small panel of cancer cell lines.</stuff_relation>

```

Fig. 3. Examples of sentences annotation.

### B. Scientific Name Entity Determination

1) *Chemical Name Entity Determination*: There are 2 steps involved: the translation of numeric representation of chemicals and the identification of chemical name entities.

a) *Translating Numeric Representation of Chemicals*: The following observed rule are used to translate the numeric representations of chemicals.

Rule:

if((Cname is the first occurrence)  $\wedge$  (Cname(num)) then  
num is the numeric representation of Cname

Where Cname = chemical name on a scientific document  
num = the integer

b) *Chemical Name Entity Identification*: The results from a) coupled with the adjacent surrounding word are compared to the NBIC-PubChem.

2) *Genus-species Name Entity Determination*: Using NBIC-Taxonomy database solves the Genus-species Name Entity.

### C. Stuff Relation Learning

In the learning step, b (where  $b \in B$  and  $B$  is obtained from the corpus behavior studied) is used to anchor sentences. Sentences where the left hand side of the anchor contains a (where  $a \in A$ ) or c (where  $c \in C$ ) and the right hand side of the anchor contains c or a are collected. Then the frequency of a, and c with class "yes" and class "no" are determined for each sentences (see TABLE III.)

TABLE III. FREQUENCY OF A AND C WITH CLASS "YES"/"NO"

A	Class=yes	Class=no
3,10-dihydroxy-9-methoxypterocarpan	0.05882353	0.03571429
ethyl acetate	0.05882353	0.17857143
dibenzocycloheptene	0.3921569	0.03571429
...	...	...
C	Class=yes	Class=no
Dendrolobium	0.06976744	0.05
Platymiscium	0.34883721	0.25
Dalbergia	0.06976744	0.25
...	...	...

### D. Stuff Relation Extraction

In order to start the Stuff relation extraction process, Naive Bayes Classifier[9] shown in equation (1) is applied. The class "yes" means Stuff relation, as shown in Fig. 4.

$$\begin{aligned}
 \text{Stuff RelationClass} &= \arg \max_{\text{class} \in \text{Class}} P(\text{class} | a, c) \\
 &= \arg \max_{\text{class} \in \text{Class}} P(a | \text{class})P(c | \text{class})P(\text{class})
 \end{aligned} \tag{1}$$

where  $\text{Class} = \{ "yes", "no" \}$

$a \in A$  ( $A$  is a Stuff word set)

$c \in C$  ( $C$  is an object word set)

```

L is a list of sentence.
W is a set of word in each sentence.
WS is a sequence of n words where n = 1, 2, 3, ...
C is natural-product-compounds concept set.
S is the natural-sources concept set.
A is the marker set.
1 i ← 0, R ← φ
2 while (i ≤ length[L]) do
3   Begin
4     j ← 0, len ← length[Wi]
5     while (wij ∈ A) do
6       for (k ← 0 to j)
7         if (wsik ∈ C) then
8           for (m ← j+1 to len)
9             if ((wsim ∈ S) ∧ (StuffRelationClass
10              = "yes")) then
11               R = R ∪ {(wsik, wsim)}
12             end if
13           next
14         else if (wsik ∈ S) then
15           for (m ← j+1 to len)
16             if ((wsim ∈ C) ∧ (StuffRelationClass
17              = "yes")) then
18               R = R ∪ {(wsim, wsik)}
19             end if
20           next
21         endif
22       next
23     End
24   return R

```

Fig. 4. Stuff Relation Extraction Algorithm.

### V. EVALUATION AND CONCLUSION

The English corpora of the technical documents in chemistry domain are used to evaluate the proposed stuff relation extraction algorithm consisting of about 5,000 sentences. The evaluation of the Stuff Relation extraction performance in research is expressed in terms of the precision and the recall as shown below, where R is the stuff relation:

$$\text{Precision} = \frac{\# \text{ of samples correctly extracted as R}}{\# \text{ of all samples output as being R}} \quad (2)$$

$$\text{Recall} = \frac{\# \text{ of samples correctly extracted as R}}{\# \text{ of all samples holding the target relation R}} \quad (3)$$

The results of precision and recall are evaluated by three expert judgments with max win voting. The precision of the extracted Stuff relation is 87% and 61% recall. The reason of the recall limited to 61% is misplaced-B problem where B is before or after A and C. These problems are a subject of further studies. Finally, these extracted Stuff relations are beneficial for natural product chemical ontology construction (see Fig. 5.) which will bring significant benefits to the cosmetics, pharmaceuticals and other chemicals industries.

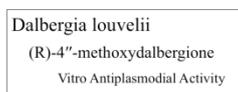


Fig. 5. Example of Natural Product Chemical Ontology.

#### REFERENCES

- [1] Bernard Munos, "Lessons from 60 years of pharmaceutical innovation," *Nature Reviews in Drug Discovery*, vol. 8, pp. 959-968, December 2009.
- [2] Yusuf Chisti, "Biodiesel from microalgae," *Biotechnology Advances*, vol. 25, pp. 294-306, February 2007.
- [3] Gail Taylor, "Biofuels and the biorefinery concept," *Energy Policy*, vol. 36, pp. 4406-4409, 2008.
- [4] Morton E. Winston, Roger Chaffin, Douglas Herrmann, "A taxonomy of part-whole relations," *Cognitive Science*, vol. 11, pp. 417-444, 1987.
- [5] Roxana Girju, Adriana Badulescu, Dan Moldovan, "Learning semantic constraints for the automatic discovery of part-whole relations," In *Proceedings of HLT/NAACL-03*, vol. 1, pp. 1-8, 2003.
- [6] Patrick Pantel, Marco Pennacchiotti, "Espresso: Leveraging generic patterns for automatically harvesting semantic relations," In *Proceedings of COLING/ACL-06*, pp. 113-120, 2006.
- [7] Katrin Fundel, "Text Mining and Gene Expression Analysis Towards Combined Interpretation of High Throughput Data," München, 2007
- [8] Gregory Ichneumon Brown, "An Error Analysis of Relation Extraction in Social Media Documents," In *Proceedings of ACL-HLT-49*, pp. 64-68, June 2011.
- [9] Tom Mitchell, *Machine Learning*. Singapore: The McGraw Hill Companies Inc. and MIT Press, 1997.