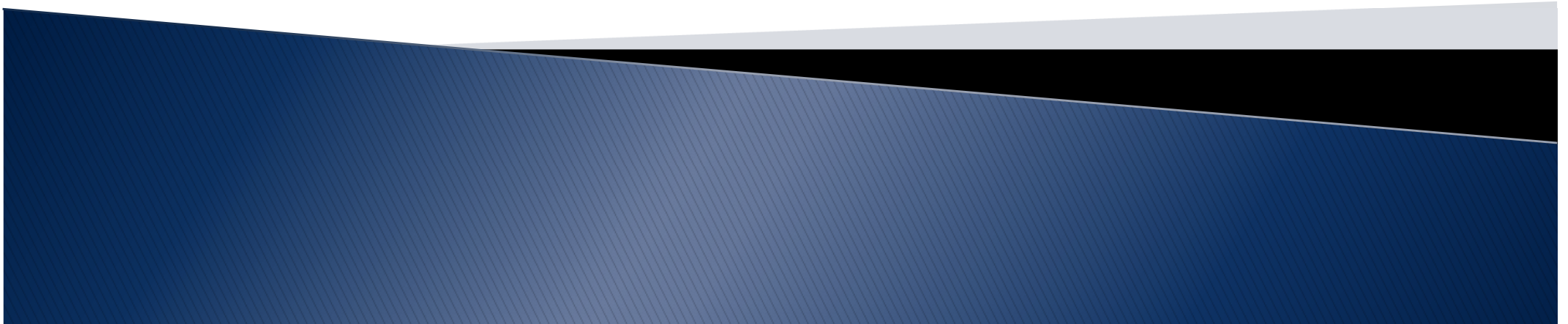


# Kobe University Brussels Workshop 7 November 2013

## The Oxford Corpus of Old Japanese

Bjarke Frellesvig  
University of Oxford



# Oxford Corpus of Old Japanese

A comprehensive digital, annotated corpus of all extant texts in Japanese from the Old Japanese period (mainly 8<sup>th</sup> century AD, Nara period)

url: <http://vsarpj.orinst.ox.ac.uk/corpus/>



## The Oxford Corpus of Old Japanese オックスフォード上代日本語コーパス

[Home](#)[Texts](#)[Display  
conventions](#)[Tagging  
conventions](#)[Searching the  
OCOJ](#)[Work in progress](#)[Events and news](#)[People and  
Institutions](#)[Funding bodies](#)[Links](#)[Contact](#)

### Welcome to the Oxford Corpus of Old Japanese!

The Oxford Corpus of Old Japanese (abbreviated OCOJ) is a long-term research project which aims to develop a comprehensive annotated digital corpus of all extant texts in Japanese from the Old Japanese period. Old Japanese is the earliest attested stage of the Japanese language, largely the Japanese language of the Asuka and Nara periods of Japanese history (7th and 8th century AD). This is the formative literate period upon which the development of Japanese civilization is based, and these texts are of paramount importance for the study and understanding of the origins and development of civilization of Japan, including language, writing, literature, religion, history, and culture.

The OCOJ will contain

1. **Texts:** The corpus will contain all extant texts in Japanese from the Old Japanese period, presented in original script and phonemic transcription. See [here](#) for the texts and [here](#) for display conventions.
2. **Annotation:** A large amount of information about the texts will be encoded and made searchable. This will include linguistic information (orthographic, phonological, morphological, syntactic, semantic, and lexical information), as well as literary, biographical, historical, geographical and other information. The digital format makes it possible to add information of any kind continuously. See [here](#) for tagging conventions.

# Funding and collaboration

## External funding

Arts and Humanities Research Council, UK

British Academy

JSPS/Tsukuba University

## Collaborating institutions

National Institute for Japanese Language and  
Linguistics, Tokyo (NINJAL 国立国語研究所)

University of York

# OCOJ Content

**Texts**

**Translations**

**Dictionary**

**Annotation**

# Texts

The corpus will contain all extant texts in Japanese from the Old Japanese period, presented in **original script** and **phonemic transcription**.



# Texts in the OCOJ

## Poetic texts

<i>Kojiki kayō</i> (古事記歌謡; 712)	(112 poems; 2527 words)
<i>Nihon shoki kayō</i> (日本書紀歌謡; 720)	(133 poems; 2444 words)
<i>Fudoki kayō</i> (風土記歌謡; 730s)	(20 poems; 271 words)
<i>Bussukoseki-ka</i> (仏足石歌; after 753)	(21 poems; 337 words)
<i>Man'yōshū</i> (万葉集; after 759)	(4685 poems; 83706 words)
<i>Shoku nihongi kayō</i> (続日本紀歌謡; 797)	(8 poems; 134 words)
<i>Jōgū shōtoku hōō teisetsu</i> (上宮聖徳法王帝説)	(4 poems; 60 words)

## ‘Prose’ texts

<i>Shoku nihongi Senmyō</i> (続日本紀宣命; 797)	(approx. 14,000 words)
<i>Engishiki Norito</i> (延喜式祝詞)	(approx. 6,500 words)

# Texts in the OCOJ

## Poetic texts

<i>Kojiki kayō</i> (古事記歌謡; 712)	(112 poems; 2527 words)
<i>Nihon shoki kayō</i> (日本書紀歌謡; 720)	(133 poems; 2444 words)
<i>Fudoki kayō</i> (風土記歌謡; 730s)	(20 poems; 271 words)
<i>Bussukoseki-ka</i> (仏足石歌; after 753)	(21 poems; 337 words)
<i>Man'yōshū</i> (万葉集; after 759)	(4685 poems; 83706 words)
<i>Shoku nihongi kayō</i> (続日本紀歌謡; 797)	(8 poems; 134 words)
<i>Jōgū shōtoku hōō teisetsu</i> (上宮聖徳法王帝説)	(4 poems; 60 words)

## ‘Prose’ texts

<i>Shoku nihongi Senmyō</i> (続日本紀宣命; 797)	(approx. 14,000 words)
<i>Engishiki Norito</i> (延喜式祝詞)	(approx. 6,500 words)
<b><i>Kojiki</i> (古事記; 712)</b>	<b>???</b>



# Texts

## **Online digital text**

website: <http://vsarpj.orinst.ox.ac.uk/corpus/>

古之

inisipye no

七賢

nana no sakasi-ki

人等毛

pito-domo *mo*

欲為物者

pori se-si mono pa

酒西有良師

sake *ni si* aru *rasi***MYS.3.341**

賢跡

sakasi-mito

物言從者

mono-ipu yworì pa

酒飲而

sake nomite

醉哭為師

wepi-naki suru *si*

益有良之

masari-taru *rasi***MYS.3.342**

將言為便

ipa-mu subye

將為便不知

se-mu subye sira-zu

極

kipamarite

貴物者

taputwo-ki mono pa

酒西有良之

sake *ni si* aru *rasi*

# Translations

The texts are being supplied with **translations** into English.

# Dictionary

**A bilingual Old Japanese – English dictionary** is being developed alongside and as an integrated part of the corpus. The dictionary part of the OCOJ is linked to the texts, making cross-reference in both directions possible.

# Annotation

A large amount of information about the texts is being encoded and made searchable.

This will include **linguistic** information (orthographic, phonological, morphological, syntactic, semantic, and lexical information), as well as **literary, biographical, historical, geographical** and other information. The digital format makes it possible to add information of any kind continuously.

# Annotation

XML mark-up following the internationally recognized standards of the **T**ext **E**ncoding **I**nitiative (TEI)

Manual mark-up

# Annotation

Linguistic information



# Annotation

Writing 原文の表記法

Part-of-speech 品詞

Lexeme and morpheme UID 語彙素と形態素のID番号

Morphology, inflection 活用

Syntax シンタクス (要素・構造)

Syntactic analysis marks and delimits **sentences, clauses,**  
and **phrases.**

# *Man'yōshū* (万葉集) 3.341

賢跡  
物言從者  
酒飲而  
醉哭為師  
益有良之

# *Man'yōshū* (万葉集) 3.341

賢跡

sakasi-mito サカシミト

物言従者

mono-ipu ywori pa モノイフヨリハ

酒飲而

sake nomite サケノミテ

酔哭為師

wepi-naki suru si エヒナキスルシ

益有良之

masari-taru rasi マサリタルラシ

“It should be better to drink saké and weep drunkenly  
than talking in a clever fashion”

```

<s>
  <phr>
    <cl>
      <cl>
        <w>
          <w type="adjective">
            <c type="logo">sakasi</c>
          </w>
          <m inflection="gerund" lemma="L000034" lemmaRef="36226" type="adjcop">
            <c type="logo">mi</c>
            <c type="phon">to</c>
          </m>
        </w>
      </cl>
      <lb corresp="#MYS.3.341-orig_1" xml:id="MYS.3.341-trans_1"/>
      <w lemma="L031771a">
        <w lemma="L050042">
          <c type="logo">mono</c>
        </w>
        <w inflection="adnconc" lemma="L030199a" lemmaRef="1571" type="verb">
          <c type="logo">ipu</c>
        </w>
      </w>
    </cl>
    <w function="abl" lemma="L000540" subtype="case" type="particle">
      <c type="logo">ywori</c>
    </w>
    <w lemma="L000522" lemmaRef="29321" subtype="top" type="particle">
      <c type="logo">pa</c>
    </w>
  </phr>
  <lb corresp="#MYS.3.341-orig_2" xml:id="MYS.3.341-trans_2"/>

```

# *Man'yōshū* (万葉集) 3.341

sakasi<sub>(adjective)</sub>-mito<sub>(adjcop, gerund)</sub>

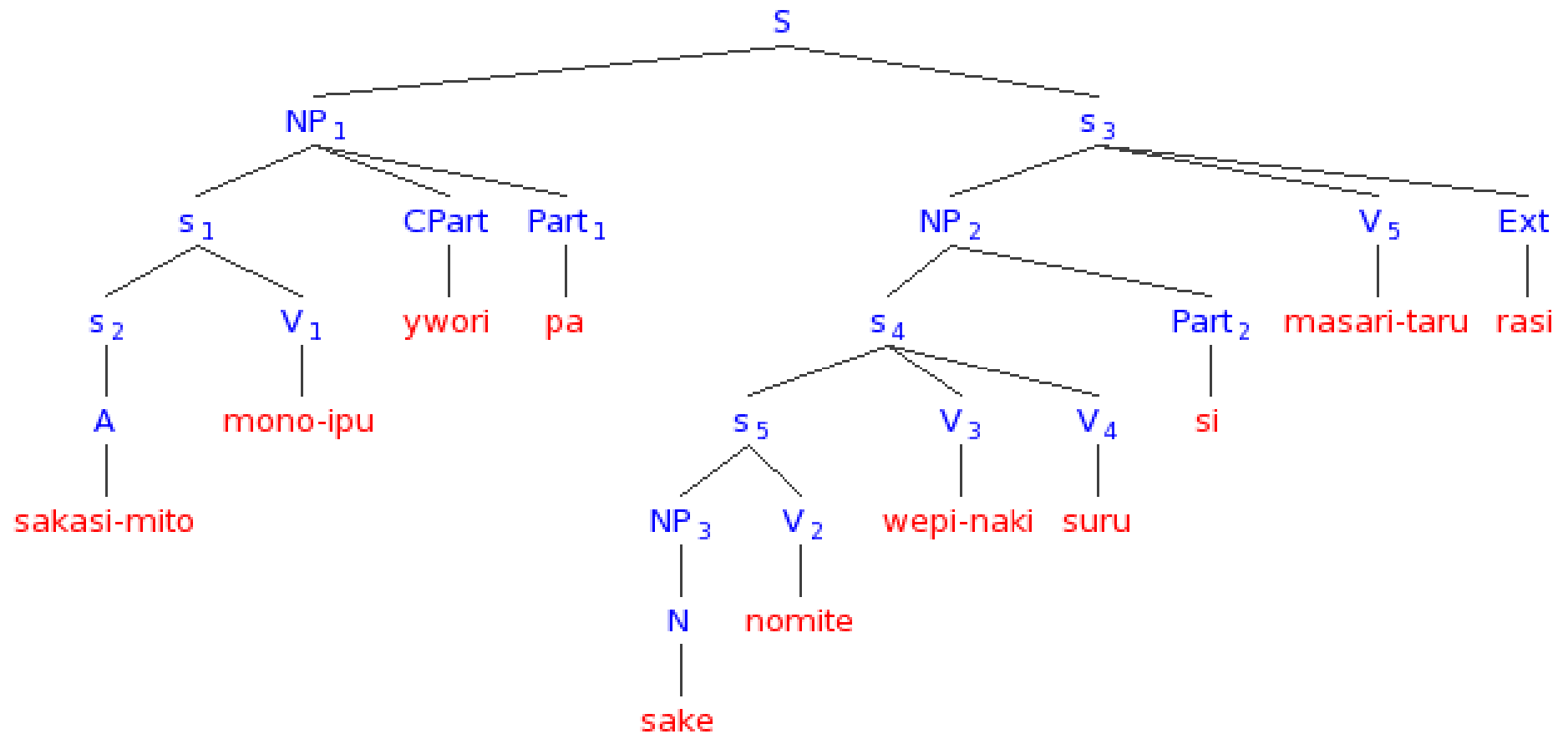
mono<sub>(noun)</sub>-ipu<sub>(verb, adnconc)</sub> ywori<sub>(abl case particle)</sub> pa<sub>(top particle)</sub>

sake<sub>(noun)</sub> nomite<sub>(verb, gerund)</sub>

wepi<sub>(verb, stem)</sub>-naki<sub>(verb, infinitive)</sub> suru<sub>(verb, adnominal)</sub> si<sub>(res particle)</sub>

masari<sub>(verb, stem)</sub>-taru<sub>(auxiliary, stative, adnominal)</sub> rasi<sub>(extension, conclusive)</sub>

# *Man'yōshū* (万葉集) 3.341



# 検索 Searching

The key to using an annotated text corpus is searchability.



# 検索 Searching

Lemmatization allows us to search the corpus for words regardless of their inflected form by doing a simple KWIC (KeyWord In Context) search:

nomu	drink.conclusive	‘I drink’
nomite	drink.gerund	‘I am drinking and’
noma-nu	drink-negative.adnominal	‘I who am not drinking’

No.	Source	Before	Search lemma	After
1	...	...	...	...
2	...	...	...	...
3	...	...	...	...
4	...	...	...	...
5	...	...	...	...
6	...	...	...	...
7	...	...	...	...
8	...	...	...	...
9	...	...	...	...
10	...	...	...	...
11	...	...	...	...
12	...	...	...	...
13	...	...	...	...
14	...	...	...	...
15	...	...	...	...
16	...	...	...	...
17	...	...	...	...
18	...	...	...	...
19	...	...	...	...
20	...	...	...	...
21	...	...	...	...
22	...	...	...	...
23	...	...	...	...
24	...	...	...	...
25	...	...	...	...
26	...	...	...	...
27	...	...	...	...
28	...	...	...	...
29	...	...	...	...
30	...	...	...	...
31	...	...	...	...
32	...	...	...	...
33	...	...	...	...
34	...	...	...	...
35	...	...	...	...
36	...	...	...	...
37	...	...	...	...
38	...	...	...	...
39	...	...	...	...
40	...	...	...	...
41	...	...	...	...
42	...	...	...	...
43	...	...	...	...
44	...	...	...	...
45	...	...	...	...
46	...	...	...	...
47	...	...	...	...
48	...	...	...	...
49	...	...	...	...
50	...	...	...	...
51	...	...	...	...
52	...	...	...	...
53	...	...	...	...
54	...	...	...	...
55	...	...	...	...
56	...	...	...	...
57	...	...	...	...
58	...	...	...	...
59	...	...	...	...
60	...	...	...	...
61	...	...	...	...
62	...	...	...	...
63	...	...	...	...
64	...	...	...	...
65	...	...	...	...
66	...	...	...	...
67	...	...	...	...
68	...	...	...	...
69	...	...	...	...
70	...	...	...	...
71	...	...	...	...
72	...	...	...	...
73	...	...	...	...
74	...	...	...	...
75	...	...	...	...
76	...	...	...	...
77	...	...	...	...
78	...	...	...	...
79	...	...	...	...
80	...	...	...	...
81	...	...	...	...
82	...	...	...	...
83	...	...	...	...
84	...	...	...	...
85	...	...	...	...
86	...	...	...	...
87	...	...	...	...
88	...	...	...	...
89	...	...	...	...
90	...	...	...	...
91	...	...	...	...
92	...	...	...	...
93	...	...	...	...
94	...	...	...	...
95	...	...	...	...
96	...	...	...	...
97	...	...	...	...
98	...	...	...	...
99	...	...	...	...
100	...	...	...	...

[Updated: 2012-Dec-0]

# Searching

More importantly,  
any combination of information  
contained in the corpus can be  
searched.

# Annotation

‘Other’ information

# Annotation

The annotation is being expanded to include amongst other information such as the following, both *in* and *about* the texts

- literary
- biographical
- historical
- geographical
- other information

What can we do for the next  
generation of Japanese studies  
scholars?

What can we do for the next  
generation of Japanese studies  
scholars?

And what can they do for us?



ご清聴ありがとうございました