

# *Geeta: Gold Standard Annotated Data, Analysis and its Application*

Preeti Shukla  
Amba Kulkarni  
Devanand Shukl

*Department of Sanskrit Studies,  
University of Hyderabad,  
Hyderabad*

20/12/2013



# Outline

- 1 Gold standard data and its advantages
- 2 Levels of Tagging
- 3 Methodology
- 4 Quantitative Analysis of BhG
- 5 Utility and Conclusion



# Gold standard data

When a corpus is manually annotated at various levels, such as Tree bank annotation, Discourse level annotation, etc., then it serves as a **gold standard data** for evaluation and comparison of various tools.

This has the following advantages –



# Advantages

- 1 One can use the gold standard annotated data as an input for the evaluation of various modules. This avoids the cascading effects and one can measure the absolute performance of various tools.
- 2 To use data-driven models to tune the machine for better performance in a chosen domain.  
For eg., Penn Tree bank



# Advantages

- 1 One can use the gold standard annotated data as an input for the evaluation of various modules. This avoids the cascading effects and one can measure the absolute performance of various tools.
- 2 To use data-driven models to tune the machine for better performance in a chosen domain.  
For eg., Penn Tree bank



# Śrīmad Bhagvad Gītā

We chose Śrīmad Bhagvad Gītā (BhG for short) consisting of 18 chapters and 700 verses (ślokas) for developing a gold standard data for Sanskrit text mainly because:

- It is an important text which summarizes the Upanishadic teachings and is commented upon and interpreted by various schools of Indian philosophies.  
Thus annotators in doubt can always refer to these commentaries for correct annotation.



- This scripture being coherent and complete in itself, can be used for higher level analysis such as discourse analysis, topic identification, anaphora resolution, and so on.



- Further, this also being part of the Mahābhārata (section 25 to 42 of the Bhiṣma parva), later on if necessary, it can be used as an initial training data for boot-strapping for automatic annotation of complete critical edition of Mahābhārata (with around hundred thousand verses).



Sanskrit requires the following levels of tagging viz.,

- 1 Annotation of Sandhi
- 2 Tagging of Compounds
- 3 Morphological Analysis
- 4 Tagging of Sentential relations
- 5 Marking the Prose order



Sanskrit requires the following levels of tagging viz.,

- 1 Annotation of Sandhi
- 2 Tagging of Compounds
- 3 Morphological Analysis
- 4 Tagging of Sentential relations
- 5 Marking the Prose order



Sanskrit requires the following levels of tagging viz.,

- 1 Annotation of Sandhi
- 2 Tagging of Compounds
- 3 Morphological Analysis
- 4 Tagging of Sentential relations
- 5 Marking the Prose order



Sanskrit requires the following levels of tagging viz.,

- 1 Annotation of Sandhi
- 2 Tagging of Compounds
- 3 Morphological Analysis
- 4 Tagging of Sentential relations
- 5 Marking the Prose order



Sanskrit requires the following levels of tagging viz.,

- 1 Annotation of Sandhi
- 2 Tagging of Compounds
- 3 Morphological Analysis
- 4 Tagging of Sentential relations
- 5 Marking the Prose order



# Annotation of Sandhi

The first level of tagging needed for Sanskrit is the marking of word boundaries undoing the sandhi (euphonic changes). We split two types of sandhis –

- 1 sandhi between two words which is indicated by a ‘+’ sign
- 2 sandhi between the components of a compound is indicated by a ‘-’ sign



# Annotation of Sandhi

The first level of tagging needed for Sanskrit is the marking of word boundaries undoing the sandhi (euphonic changes). We split two types of sandhis –

- 1 sandhi between two words which is indicated by a ‘+’ sign
- 2 sandhi between the components of a compound is indicated by a ‘-’ sign



# Example

BhG 2.40 –

**nehābhikramanāśo'asti pratyavāyo na vidyate|  
svalpamapyasya dharmasya trāyate mahato bhayāt||**

Eng: In this path there is no loss of effort, nor is there any adverse result. Even a little practice of this discipline protects one from great fear (of birth and death).



## Verse tokenized as

*na+iha+abhikrama-nāśah+asti pratyavāyah+na vidyate|  
svalpam+api+asya dharmasya trāyate mahataḥ+bhayāt||*



# Tagging of Compounds

Sanskrit compounds are broadly classified as follows:

- ① endo-centric / adverbial (tatpuruṣaḥ / avyayībhāvaḥ)
- ② exo-centric (bahuvrīhiḥ)
- ③ copulative (karmadhārayaḥ)
- ④ conjunctive (dvandvaḥ)

These compounds with an exception of dvandva (conjunctive) and bahuvrīhi (exocentric) are binary.



# Tagging of Compounds

Sanskrit compounds are broadly classified as follows:

- 1 endo-centric / adverbial (tatpuruṣaḥ / avyayībhāvaḥ)
- 2 exo-centric (bahuvrīhiḥ)
- 3 copulative (karmadhārayaḥ)
- 4 conjunctive (dvandvaḥ)

These compounds with an exception of dvandva (conjunctive) and bahuvrīhi (exocentric) are binary.



# Tagging of Compounds

Sanskrit compounds are broadly classified as follows:

- ① endo-centric / adverbial (tatpuruṣaḥ / avyayībhāvaḥ)
- ② exo-centric (bahuvrīhiḥ)
- ③ copulative (karmadhārayaḥ)
- ④ conjunctive (dvandvaḥ)

These compounds with an exception of dvandva (conjunctive) and bahuvrīhi (exocentric) are binary.



# Tagging of Compounds

Sanskrit compounds are broadly classified as follows:

- ① endo-centric / adverbial (tatpuruṣaḥ / avyayībhāvaḥ)
- ② exo-centric (bahuvrīhiḥ)
- ③ copulative (karmadhārayaḥ)
- ④ conjunctive (dvandvaḥ)

These compounds with an exception of dvandva (conjunctive) and bahuvrīhi (exocentric) are binary.



The meaning of a compound is decided by the way the components are combined together.

A compound with 3 components a-b-c may be combined in two different ways viz.,

- $\langle\langle a-b \rangle - c \rangle$
- $\langle a - \langle b-c \rangle \rangle$



# Example

The compound word **dehāntaraprāptiḥ** in BhG 2.13 is combined as  
 <<**a-b**>-c>

**anyaḥ dehaḥ** – **dehāntaram**

**Gloss:** transference\_of\_the\_body

**dehāntarasya prāptiḥ** – **dehāntaraprāptiḥ**

**Gloss:** attainment\_of\_transference\_of\_the\_body



# Example

The compound word **dehāntaraprāptiḥ** in BhG 2.13 is combined as  
 <<a-b>-c>

**anyaḥ dehaḥ** – **dehāntaram**

**Gloss:** transference\_of\_the\_body

**dehāntarasya prāptiḥ** – **dehāntaraprāptiḥ**

**Gloss:** attainment\_of\_transference\_of\_the\_body



# Various stages

Kumar et al., (2010) describe various stages involved in the analysis of a compound which also form the natural modules of a compound processor –

- ① Segmentation (samāsapadacchedadh)
- ② Constituency Parsing (samāsapadānvayaḥ)
- ③ Compound Type Identification (samastapadaparicāyakaḥ)
- ④ Paraphrasing (vighraha-vākyam)

A hierarchical tagset of 55 tags has been designed to tag the Sanskrit compounds.



# Various stages

Kumar et al., (2010) describe various stages involved in the analysis of a compound which also form the natural modules of a compound processor –

- ① Segmentation (samāsapadacchedadh)
- ② Constituency Parsing (samāsapadānvayaḥ)
- ③ Compound Type Identification (samastapadapariçāyakaḥ)
- ④ Paraphrasing (vighraha-vākyam)

A hierarchical tagset of 55 tags has been designed to tag the Sanskrit compounds.



# Various stages

Kumar et al., (2010) describe various stages involved in the analysis of a compound which also form the natural modules of a compound processor –

- ① Segmentation (samāsapadacchedadh)
- ② Constituency Parsing (samāsapadānvayaḥ)
- ③ Compound Type Identification (samastapadaparicāyakaḥ)
- ④ Paraphrasing (vighraha-vākyam)

A hierarchical tagset of 55 tags has been designed to tag the Sanskrit compounds.



# Various stages

Kumar et al., (2010) describe various stages involved in the analysis of a compound which also form the natural modules of a compound processor –

- ① Segmentation (samāsapadacchedadh)
- ② Constituency Parsing (samāsapadānvayaḥ)
- ③ Compound Type Identification (samastapadaparicāyakaḥ)
- ④ Paraphrasing (vighraha-vākyam)

A hierarchical tagset of 55 tags has been designed to tag the Sanskrit compounds.



# Morphological Analysis

At this stage tagging of both the inflectional as well as derivational information is needed. e.g.,

The morph of the verb 'asti' (to be) is

V

as2{active;laṭ;1st;sg;parasmaipadī;asaz;adādiḥ}



# Tagging of Sentential relations

In order to interpret the meaning of a sentence, various relations among words are necessary.

Of around **90 relations** classified by Ramakrishnamacharyulu (2009), only around **35 relations** based on syntactico-semantic information available in a sentence are considered for automatic tagging.



# Marking the Prose order

For understanding a Sanskrit text in verse style, two different methods have been followed in Indian education system viz.,

- **Daṇḍānvaya** (also known as anvayamukhī)

The teacher arranges all the words in prose order for easy understanding of a verse.

This approach assumes that if a user is given the 'default prose order' of the sentence, he 'understands' its meanings.



# Example

BhG 4.8 –

**paritrāṇāya sādḥūnām vināshāya ca duṣkṛtām |  
dharmasaṃsthāpanārthāya sambhavāmi yuge yuge ||**

Anvaya: *aham sādḥūnām paritrāṇāya duṣkṛtām vināshāya  
dharmasaṃsthāpanārthāya ca yuge yuge sambhavāmi*

Eng: I appear from time to time for protecting the good, for transforming the evil-minded, and for establishing world order (Dharma).



- **Khaṇḍānvaya** (also known as kathambhūtinī)

The teacher gives the basic skeleton of a sentence and fills in other details by asking questions which are centered around the heads seeking their various modifiers.

This approach is close to parsing a sentence showing various dependency relations.



# Canonical Form

The default word order or the 'canonical form' is governed roughly by the following verse:

samāsacakram kā.verse 10

*viśeṣaṇam puraskṛtya viśeṣyam tad-lakṣaṇam|  
karṭṛ-karma-kriyā-yuktam etad anvaya-lakṣaṇam||*

gloss: Starting with the adjectives, targeting the headword, in the order of karṭṛ-karma-kriyā (subject-object-verb) gives an anvaya.



# Example – BhG 2.40

```

<l>
<seg type="pāda">
<euphonic-word no="1"> nehābhikramanāśo'sti
<word no="1" prose word order_no="3"> na
<mo_anal> na{ind} </mo_anal>
<syntactic_rel> mod 4 </syntactic_rel>
</word>
<word no="2" prose word order_no = "1"> iha
<mo_anal> iha{ind} </mo_anal>
<syntactic_rel> loc 4 </syntactic_rel>
</word>

```



```

<word no="3" prose word order_no="2"> abhikramanāśaḥ
<compound label="T6">
<component no="1"> abhikrama </component>
<component no="2"> nāśaḥ
<mo_anal> nāśa{masc}{nom;sg} </mo_anal>
</component>
</compound>
<syntactic_rel> subj 4 </syntactic_rel>
</word>
<word no="4" prose word order_no="4" > asti
<mo_anal>
as2{active;lat;1st;sg;parasmaipadī;asaz;adādiḥ} </mo_anal>
</word>
</euphonic-word>
</seg>
</l>

```



The process for semi-automatic tagging of BhG is as follows:

- The verse form is converted into prose form.



- Initially the sandhi and compound in the verse are segmented manually, following the guidelines developed by the SHMT consortium <sup>1</sup>. Then each compound is tagged for its type, along with the complete constituency mark-up.

---

<sup>1</sup>This is the Consortium of 7 institutes, for 'Development of Sanskrit-Hindi Machine Translation System (sampark)' funded by DIT, Govt. of India



- The segmented words are run in the anusāraka interface <sup>2</sup> for obtaining the multiple morph analysis. The output generated as an xml file, is then manually pruned for choosing the correct morph analysis in the context.

---

<sup>2</sup><http://sanskrit.uohyd.ernet.in/scl>



- The syntactico-semantic relations are tagged manually, following the guidelines developed by the SHMT consortium <sup>3</sup>.

---

<sup>3</sup><http://sanskrit.uohyd.ernet.in/scl/Corpus/TaggingGuidelines/kaaraka-tagging-guidelines>



- Hindi and English glosses for each word are given manually. For this we followed *Geeta Press*<sup>4</sup>.

---

<sup>4</sup>Śrīmad Bhagvad Gītā, Geeta Press, Gorakhpur, India, reprint 2007



# Compound Distribution

Compound-type	Freq
Endocentric	994
Exocentric	390
Copulative	163
Conjunctive	144

◀ Refer DSP



# Morphological statistics:

multiple morph count	words	multiple morph count	words
1	5280	8	22
2	1570	9	28
3	1082	10	15
4	349	11	10
5	292	12	7
6	99	13	3
7	121	14	6
		Total words	8884
		Average	1.90



## Case Statistics in Geeta

case	singular	dual	plural
nom.	2463	31	613
acc.	1349	20	251
instr.	266	0	94
dat.	57	0	5
abl.	116	0	12
gen.	335	24	179
loc.	273	3	92
voc.	251	1	0





# syntactico-semantic relations in Geeta

relation	freq	relation	freq
viśeṣaṇa (adjective)	1277	kartā (subject)	1256
conjunctive	1155	karma (object)	924
predicative adj	401	adhikaraṇa (locative)	358
ṣaṣṭhi (genitive)	357	negation	242
emphatic	275	sambodhana (vocative)	237
precedence	194	adverb	194
karaṇa (instrument)	130	karma-	113
hetu (causal)	96	sāmānidhikaraṇa	
apādāna (source)	75	co-relative	82
prayojana (purpose)	52	sarvanāma (pronouns)	
sampradāna (dative)	15	vākya-karma	64
		simultaneity	50
		disjunction	10



# Tense-Mood distribution in Geeta

Sanskrit has 10 lakāras which represent Tense-Modality.

lakāra	freq
laṭ (Present)	355
loṭ (Imperative)	72
lṛṭ (Second future)	42
vidhiliṅ (Potential)	40
laṅ (Imperfect)	26
liṭ (Perfect)	22
luṭ (First future)	16
luṅ (Aorist)	9
āśīrliṅ (Optative)	6
lṛṅ (Conditional)	1



# Utility

There are three important usages of this gold data viz.,

- 1 for NLP applications
- 2 as linguistic inputs for developing a domain specific primer for learning BhG
- 3 with the suitable interface, a self-learning / reading tool for BhG.



# Utility

There are three important usages of this gold data viz.,

- 1 for NLP applications
- 2 as linguistic inputs for developing a domain specific primer for learning BhG
- 3 with the suitable interface, a self-learning / reading tool for BhG.



# Utility

There are three important usages of this gold data viz.,

- 1 for NLP applications
- 2 as linguistic inputs for developing a domain specific primer for learning BhG
- 3 with the suitable interface, a self-learning / reading tool for BhG.



# NLP applications

- serves as a gold standard for evaluation of various Sanskrit tools.
- for developers of NLP tools.
- in prioritizing the solutions in the case of ambiguities.



# NLP applications

- serves as a gold standard for evaluation of various Sanskrit tools.
- for developers of NLP tools.
- in prioritizing the solutions in the case of ambiguities.



# NLP applications

- serves as a gold standard for evaluation of various Sanskrit tools.
- for developers of NLP tools.
- in prioritizing the solutions in the case of ambiguities.



# NLP applications

- serves as a gold standard for evaluation of various Sanskrit tools.
- for developers of NLP tools.
- in prioritizing the solutions in the case of ambiguities.



# Domain Specific Primer

- The quantitative analysis can help a teacher to decide which aspect of Sanskrit Grammar is more relevant for the study of BhG as follows:
  - 1 postponing the teaching of dual forms to a later stage.
  - 2 deciding how many and which paradigms of noun declension to concentrate on first.
  - 3 using the analysed data of Tense-Modality to decide which lakāras to teach and which conjugation classes to concentrate on first.
  - 4 deciding on which type of compound to teach first based on the compound distribution table.

[DSP](#)[▶ Go to Case](#)[▶ Go to lakara](#)[▶ Go to CD](#)

# Domain Specific Primer

- The quantitative analysis can help a teacher to decide which aspect of Sanskrit Grammar is more relevant for the study of BhG as follows:
  - ① postponing the teaching of dual forms to a later stage.
  - ② deciding how many and which paradigms of noun declension to concentrate on first.
  - ③ using the analysed data of Tense-Modality to decide which lakāras to teach and which conjugation classes to concentrate on first.
  - ④ deciding on which type of compound to teach first based on the compound distribution table.

[DSP](#)
[▶ Go to Case](#)
[▶ Go to lakara](#)
[▶ Go to CD](#)


# Self Learning cum Reading Tool

With the help of a suitable interface such as that of anusāraka, an interested reader can have complete analysis of BhG at various levels. The interface provides –



- User controlled access to various levels of analysis (Figs. 1,2,3). The graphs showing the constituency information and the kāraka relations are generated automatically from the manually tagged data.

श्रीमद्भगवद्गीता  
गीता सुगीता कर्तव्या किमन्यद् शास्त्र विरक्तेः

सुखदुःखे समं कृत्वा लाभालाभौ जयाजयम्  
ततो बुद्ध्या युज्यस्व नैवम पापमाकर्ष्यसि ।।2.38।।

एषा तेजिता संख्ये बुद्धिसौ त्विमां श्रुत्वा  
बुद्ध्या युक्तो यया पापं कर्मबन्धन प्रहास्यसि ।।2.39।।

नेत्राभ्रमनाहोसि प्रत्यवाधो न विद्यते  
स्वल्पमप्यस्य धर्मस्य उपरो महतः प्रयात् ।।2.40।।

यवसायान्तिका बुद्धिकेव कुरुनन्दन  
बहुतासा हनन्तास्य बुद्धीय्यसायिनम् ।।2.41।।

यागिभ्यमुषिणां वाचमप्रवदन्त्यपिचरिचिः  
केवद्वरताः पार्श्वं मान्यदत्तौगि वादिनः ।।2.42।।

कामाचानः स्वर्गारा जन्मकर्मफलप्रदानं  
क्रियतिशेषद्वन्द्वतन्मोक्षैर्यदितिप्रति ।।2.43।।

भौतैस्वर्गप्रदानकनन्त्यवद्वत्तैसत्साम्  
व्यवसायान्तिका बुद्धिः शमधो न विधीयते ।।2.44।।

नैमुष्यचिष्वा वेदा निरुग्धो भवार्जुन  
निर्द्वन्द्वे निरवसरथो नियोगिक्सेम आच्यमान ।।2.45।।

1.1.A	iha	abhikrama-nāśah	na	asti	pratyavāyāb	na	vidyate	asya
1.1.E		<.abhikrama-nāśah>T6						
1.1.H	isa	karmayoga_mem	ārambhakā_arthāt	bijakā_nāśa	nahim	hai	ulatā_g	ā_dosa
1.1.I	in_this_world	endeavoring_loss						never is of_this

(abhikrama\_nāśah)[T6]

abhikrama      nāśah

dharmasya	svल्पam	api						trāyate
dharma_kā	thodā-sā	bhī_(sādhana_jan						reksā_kara_letā_hā
of_this_occupation	a_little	although						releases

Show/Hide Rows... Numbers Borders anvaya.file

Figure: 1. compound analysis of BhG 2.40



श्रीमद्भगवद्गीता  
गीता सुगीता कर्तव्या विमन्वत् शास्त्र विस्तरः

युक्तदुःखं शमं कृत्वा ज्ञानोत्तमो ज्ञयत्येव  
ततो युद्धाय युक्तस्यैव नीचम् पापमवाप्तस्यसि ॥2.38॥

एषा तेषिहिता सांख्ये बुद्धिदोषो त्स्त्रिणां श्रुत्वा  
ब्रुवाया युक्तो यथा पाथं कर्मवन्धम् प्रहासयसि ॥2.39॥

येनाधिकमनात्तोस्ति प्रत्यक्षसो न विद्यते  
स्वल्पमप्यस्य धर्मस्य जगते महताः भयात् ॥2.40॥

अवसायात्मिका बुद्धिश्चेत्क्षुरणन्तु  
ब्रह्मसाक्षात्क्षानन्तास्य बुद्धयोपवसायिणम् ॥2.41॥

वाग्मिणामुष्मिनां वाक्प्रभवन्त्यधिपस्थिताः  
केवलसरताः पाथं मान्यवदन्तीति वाचिनः ॥2.42॥  
कामादानः स्वर्गपरा जन्मकर्मफलप्रदम्  
क्रियाधिभेदबहुलान्भोगैस्वयंनिप्रति ॥2.43॥  
भोगैस्वयंप्रसक्तान्प्रत्ययाद्ब्रह्मवेत्तसाम्  
अवसायात्मिका बुद्धिः समाधी न विधीयते ॥2.44॥

केनूग्रथिष्या वेदा निरतैर्गुणो भवार्जुन  
निर्दुःको नित्यसत्सथो नियोगेषु आत्मवान् ॥2.45॥

1.1.A na abhikrama-nāśah na asti pratyavāyah na vidyate asya  
1.1.E <abhikrama-nāśah=T6  
1.1.H sa\_karmayoganūpa  
1.1.I of\_this  
trāyate  
kartā apādānam  
svalpam bhayāt  
sasthīsambandhaḥ sambandhaḥ viśeṣanam  
dharmasya mahataḥ  
viśeṣanam  
asti  
anya file

Figure: 2. kāraka analysis of BhG 2.40



1.1.A iha	abhikrama-nāśaḥ	na	asti	pratyavāyāḥ	na
1.1.D iha {avya}	abhikrama-nāśa {puṃ} {1;eka}	na {avya}	as2 {kartari;lat;pra;eka;parasmaipadi;asaṃ;adādh}	pratyavāya {puṃ} {1;eka}	na {avya}
1.1.H isa_karmayoga_mem	ārambhakā_arthāt_bijakā_nāśa	nahiṃ	hai	ulaṭā_phala_kā_doṣa	na
1.1.I in_this_world	endeavoring_loss	not	is	diminution	never
vidyate	asya	dharmasya	svalpam	api	
vid2 {kartari;lat;pra;eka;ātmanepadi;vidāṃ;divādh}	idam {puṃ} {6;eka}	dharma {puṃ} {6;eka}	svalpa {napuṃ} {1;eka}	api {avya}	
hai	isa_karmayogarūpa	dharma_kā	thoḍā-sā	bhi_{sādhana_janma-mṛtyurupa}	
is	of_this	of_this_occupation	a_little	although	
mahataḥ	bhayāt	trāyate			
mahat {napuṃ} {5;eka}	bhaya {napuṃ} {5;eka}	trai1 {kartari;lat;pra;eka;ātmanepadi;train;bhvādh}			
mahān	bhayase	rakṣā_kara_letā_hai			
of_very_great	danger	releases			

Show/Hide Rows...  Numbers  Borders

[anvya file](#)

Figure: 3. morph analysis of BhG 2.40



- Link to various dictionaries for meanings of the head words.
- Graphical display of phrase structure analysis of compounds (Fig 4).

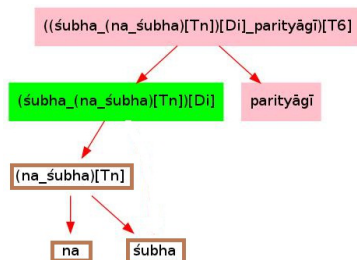


Figure: 4. Compound analysis of a word from BhG 12.17



- Graphical display of sentential analysis (Fig 5).

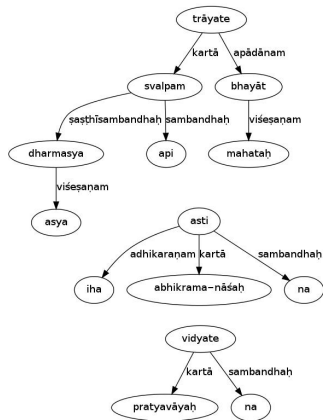


Figure: 5. Dependency graph of BhG 2.40



*This provides the user a digitized learning and understanding environment and forms a basis for the theoretical linguists and grammarians to test their theories as well.*



धन्यवाद:!!!

