

# SCCT Annotation Guidelines

Hello, dear annotator, and thank you for your time!

Here is a brief background for the task and the specific annotation guidelines. This task is a part of the DharmaBench paper and evaluation benchmark for Sanskrit and Tibetan.

## Introduction

**Task:** Classifying canonical literature (SCCT)

**Goal:** To distinguish between Tibetan allochthonous literature into the two categories of scripture and non-scripture.

**Background:** This task classifies Tibetan canonical literature (allochthonous) into two primary categories: “scriptures” and “non-scriptures.” Its primary objectives are to (a) evaluate LLM performance in canonical literature classification into these two broad categories, and (b) illuminate the evolutionary processes of works that are considered scriptures by the tradition (allochthonous texts) as well as para-canonical works through (works whose provenance is debated) linguistic analysis. The task aims at identifying distinctive features characteristic of each category, including compositional styles, syntax, vocabulary, and other linguistic markers. In particular, such an analysis should help reveal recurring stylistic patterns and compositional strategies that a scripture creator possibly employed to achieve two critical goals: canonical acceptance and popular success.

**Real-world relevancy:** By tracing and mapping these linguistic patterns, this task has the potential to provide insights into how textual authority was constructed and maintained within Indic and Tibetic literary traditions.

## Example

1. **Text:** “ཀླུ་པོའི་ཁབ་ཀྱི་གྲོང་ཁྱེར་ཆེན་པོར་བྱོན་ནས་དེ་ཀླུ་པོའི་ཁབ་ཀྱི་གྲོང་ཁྱེར་ཆེན་པོར་བྱོན་ཀྱི་ཕུང་པོའི་རི་ལ་བཞུགས་ཏེ། འདི་ལྟར་བཅོམ་ལྷན་འདས་དེ་ནི་དེ་བཞིན་གཤེགས་པ་དག་བཅོམ་པ་ཡང་དག་པར་རྫོགས་པའི་སངས་རྒྱུས་རིག་པ་དང་ཞབས་སུ་ལྷན་པ། བདེ་བར་གཤེགས་པ་འཛིག་རྟེན་མཁྱེན་པ། སྤྱི་བུ་ལ་བའི་ཁ་ལོ་སྦྱར་བ། ལྷ་ན་མེད་པ། ལྟ་དང་མི་རྣམས་ཀྱི་སྦྱོན་པ། སངས་རྒྱུས་བཅོམ་ལྷན་འདས་ཡིན་ཏེ། དེ་ཆོས་སྦྱོན་པ་ནི་ཆངས་པར་སྦྱོད་པ། ཐོག་མར་དགེ་བ། བར་དུ་དགེ་བ། ཐ་མར་དགེ་བ། དོན་བཟང་པོ། ཆོག་འབྱུང་བཟང་པོ། མ་འདྲེས་པ་ཡོངས་སུ་རྫོགས་པ། ཡོངས་སུ་དག་པ། ཡོངས་སུ་བྱང་བ་རབ་ཏུ་སྦྱོན་ཏེ་ཞེས་དགེ་སྦྱོང་གོ་ཏེ་མ་དེའི་དགེ་བར་གཤེགས་པའི་སྤྱི་ཆོག་སུ་བཅད་པ་རྒྱ་ཆེན་པོ་དེ་ལྟ་བུ་ཡང་མཛོན་པར་བྱུང་བ་ཀྱིས་བདག་ཆོང་དཔོན་བཟང་སྦྱོང་གིས་ཐོས་ནས་དེ་འདི་སྤྱི་བུ་སེམས་ཏེ། དེ་བཞིན་གཤེགས་པ་དག་བཅོམ་པ་ཡང་དག་པར་རྫོགས་པའི་སངས་རྒྱུས་དེ་ལྟ་བུ་དག་མཛོན་ན་ལེགས་པར་འབྱུང་གིས་བདག་དགེ་སྦྱོང་གོ་ཏེ་མ་དེ་བུ་

བའི་ཕྱི་ལོ་ལྷན་ཁོ་། དེ་ནས་ཁྱིམ་བདག་ཚོང་དཔོན་བཟང་སྟོང་ཁྱིམ་བདག་ལྷ་བརྒྱ་ཙམ་དང་ལྷན་ཅིག་ཏུ་བཙུགས་ལྷན་འདས་ལ་བལྟ་བའི་ཕྱི་ལོ་ལྷ་པོ་འཁྲུག་གི་

གོང་ཁྱེད་ཆེན་པོ་ནས་བྱུང་ངོ་།” Class: SCR

2. **Text:** “སངས་རྒྱལ་ལ་མ་དད་ཅེས་བྱ་བ་ནི་ཡིད་མ་ཆེས་པས་བརྟན་ཤིང་མ་དད་པའི་དོན་རྟོ། །ལྷ་གཞན་ལ་བརྟེན་ཞེས་བྱ་བ་ནི་འཇིག་རྟེན་པའི་ལྷ་ལ་བརྟེན་ཅིང་ལྷ་སྟེགས་པའི་ཕྱོགས་འཇིན་པའོ། །ང་ཡི་གས་དང་ཐུགས་བརྒྱས་ན་ཕུང་བར་འགྱུར། །ཞེས་བྱ་བ་ནི་འཆི་བ་དང་འཇིག་རྟེན་པ་རོལ་དུ་ཡང་དན་སོང་དུ་སྐྱེ་བར་འགྱུར་རོ། །བསམ་པ་ཞན་ཅིང་ཞེས་བྱ་བ་ནི་སེམས་ཁྱུམ་པའོ། །ཉམས་ཀྱིས་མི་ལྟོགས་པས་ཞེས་བྱ་བ་ནི་ཡོ་བྱད་དང་མི་ལྷན་ཞེས་བྱ་བའི་དོན་རྟོ། །དཀྱིལ་འཁོར་ཀྱན་ཞེས་བྱ་བ་ནི་གོང་དུ་སྐྱོས་པ་ཐམས་ཅད་དོ། །བལྟ་བར་མི་རུས་ན་འད་ཞེས་བྱ་བ་ནི་ཡོ་བྱད་མ་འགྱོར་ཏེ་སེམས་ཁྱུམ་ནས། དཀྱིལ་འཁོར་འཁྲིལ་འཇུག་མ་རུས་ན་ཡང་། ཞེས་བྱ་བའི་དོན་རྟོ། །གཅིག་ཙམ་ཞིག་ཀྱང་དད་པས་མཐོང་གྱུར་ན། ཞེས་བྱ་བ་ནི་ཡོ་བྱད་མེད་པའི་སྐབ་པ་པོ་དབུལ་པོས་ཀྱང་དཀྱིལ་འཁོར་དེ་དག་རྣམས་ལས་ཉམས་ཀྱིས་གང་ལྟོགས་པ་གཅིག་ནི་ངེས་པར་འདྲི་དགོས་པར་གསུངས་ཏེ། དོན་དེ་ཉིད་རབ་ཏུ་གྱུབ་པར་བྱེད་པའི་རྒྱུད་ཆེན་པོ་ལས་ཀྱང་ཇི་སྟངས་སུ།” Class: NSRC

## Annotation guidelines

1. The samples will consist of heuristically cut sections of text from canonical and paracanonical sources (ACIP & rKTs). They will be provided in a CSV. file with a pre-assigned label. Read the sample of the text carefully and determine if the label is correct.
2. If the sample contains any material that is not part of the original text, such as table of contents, editorial notes, translator’s remarks, and so forth, reject the sample.
3. If the sample is too short (i.e., shorter than two full sentences), reject the sample.
4. If there are any conversion errors in the Tibetan (due to erroneous input in the original material, for example), if the error is minor, correct it. If there are many conversion errors that cannot be resolved easily, reject the sample.

## Edge cases and common mistakes

1. In cases of ambiguity, please consult with the team and your supervisors.
  - a. If you can’t reach an agreement, discard the sample.

## Something is off?

Reach out to us.

Thank you very much!