

# Caption Label Noise Dataset (CLaN) Documentation - Data Appendix

---

First we provide an overview of the file contents included in this appendix. Then we provide details about the data collection and annotation protocol. Then, we provide information on how to setup our annotation interface. Lastly, we show tables and figures showing agreement stratified by question, option, and dataset.

## Documents

1. [veil\\_annotation\\_project.zip](#) : a compressed directory containing images from each in-the-wild dataset.
2. [annotator\\_1\\_labelstudio\\_export.json](#) : contains metadata used to populate Label Studio annotation interface and selected annotations from annotator 1.
3. [annotator\\_2\\_labelstudio\\_export.json](#) : contains metadata used to populate Label Studio annotation interface and selected annotations from annotator 2.
4. [clan-label-interface.xml](#) : contains labeling interface code which can be pasted into label-studio. Note that the alias (stored in annotations.csv) mapping is in this file.
5. [annotations.csv](#) : contains all annotations

## Data Collection and Annotation Protocol

We use three in-the-wild datasets (SBUCaps, RedCaps, Conceptual Captions) to construct this dataset. We perform weighted sampling to draw VAEs from the same distribution as WSOD training; this also represents VAEs from a diverse range of categories rather than a few frequent categories. We sample 100 VAEs, their caption, and corresponding image from each in-the-wild VL dataset.

### Annotation Schema

We annotate four types of information for each sample:

1. (Q1: Label Noise) How much of the VAE object is present (visible, partially visible, completely absent);
2. (Q2: Similar Context) If the VAE object is completely missing, whether a traditionally co-occurring context ("boat" and "water"), or semantically similar object (e.g. "cake" and "bread", "car" and "truck") is present;
3. (Q3: Visual Defects) If visible/partially visible, whether the VAE object is occluded, has key parts missing, or atypical appearance (e.g. knitted animal); and
4. (Q4: Linguistic Indicators) What linguistic cues, if any, explain why the VAE object is mentioned but absent, e.g. the caption discusses events or information beyond what the image shows ("beyond" in Table 1), describes the past ("past"), the extracted label is part of a prepositional phrase and likely to describe the setting and not objects ("on a train"), is a noun modifying another noun, is used in a non-literal way, has a different word sense (e.g. "bed" vs "river bed"), or is part of a named entity.

We compared disagreements between the two annotators after they annotated 100 samples from RedCaps and used it to calibrate responses for the remaining samples from SBUCaps and Conceptual Captions. We

indeed observe higher agreement (weighted Cohen's Kappa) for Q2-Q4 the remaining datasets in the Agreement section.

## Annotation Interface Setup

1. Download Label-Studio (requires Python version 3.8+)

```
python3 -m venv env  
source env/bin/activate  
python -m pip install label-studio
```

2. Unzip compressed image data

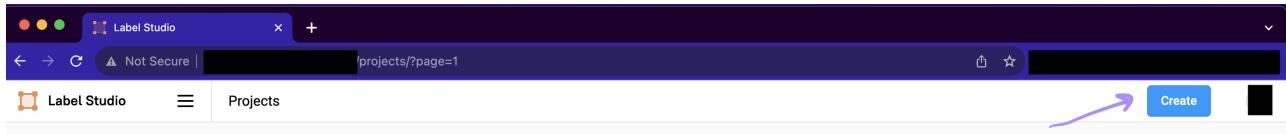
```
unzip veil_annotation_project.zip -d veil_annotation_project
```

3. Setup environment variables

```
LABEL_STUDIO_LOCAL_FILES_SERVING_ENABLED=TRUE  
LABEL_STUDIO_LOCAL_FILES_DOCUMENT_ROOT=/ # replace with absolute path to  
the parent directory of veil_annotation_project
```

4. Launch Label Studio **label-studio**

5. Create a new project.



6. Import one of the JSON files above **annotator\_1\_labelstudio\_export.json** OR **annotator\_2\_labelstudio\_export.json**. To populate the labeling interface, copy the contents of **clan-label-interface.xml**.

Create Project      Project Name      Data Import      Labeling Setup      Delete      Save

Dataset URL      Add URL      or     

**Drag & drop files here  
or click to browse**



Text      txt

Audio      wav, aiff, mp3, au, flac, m4a, ogg

Video      mpeg4/H.264 webp, webm\*

Images      jpg, png, gif, bmp, svg, webp

HTML      html, htm, xml

Time Series      csv, tsv

Common Formats      csv, tsv, txt, json

\* – Support depends on the browser

 See the documentation to [import preannotated data](#) or to sync data from a database or cloud storage.

Create Project      Project Name      Data Import      Labeling Setup      Delete      Save

**Computer Vision** >

- Natural Language Processing >
- Audio/Speech Processing >
- Conversational AI >
- Ranking & Scoring >
- Structured Data Parsing >
- Time Series Analysis >
- Videos >
- Generative AI >
- Custom template**



  
Semantic Segmentation with Polygons

  
Semantic Segmentation with Masks

  
Object Detection with Bounding Boxes

  
Keypoint Labeling

  
Multi-page document annotation

  
Inventory Tracking

 See the documentation to [contribute a template](#).

Create Project      Project Name      Data Import      Labeling Setup      Delete      Save

Code       Visual

```
1 <View>
2 </View>
3
```

**UI Preview**

Regions	History	Relations	Info
Manual	By Time ↑		
Regions not added			

7. Finally, set the target storage. This is where you want to store your annotations. Note that if you set the environment variables correctly, then all the images should load correctly. (redacted example below)

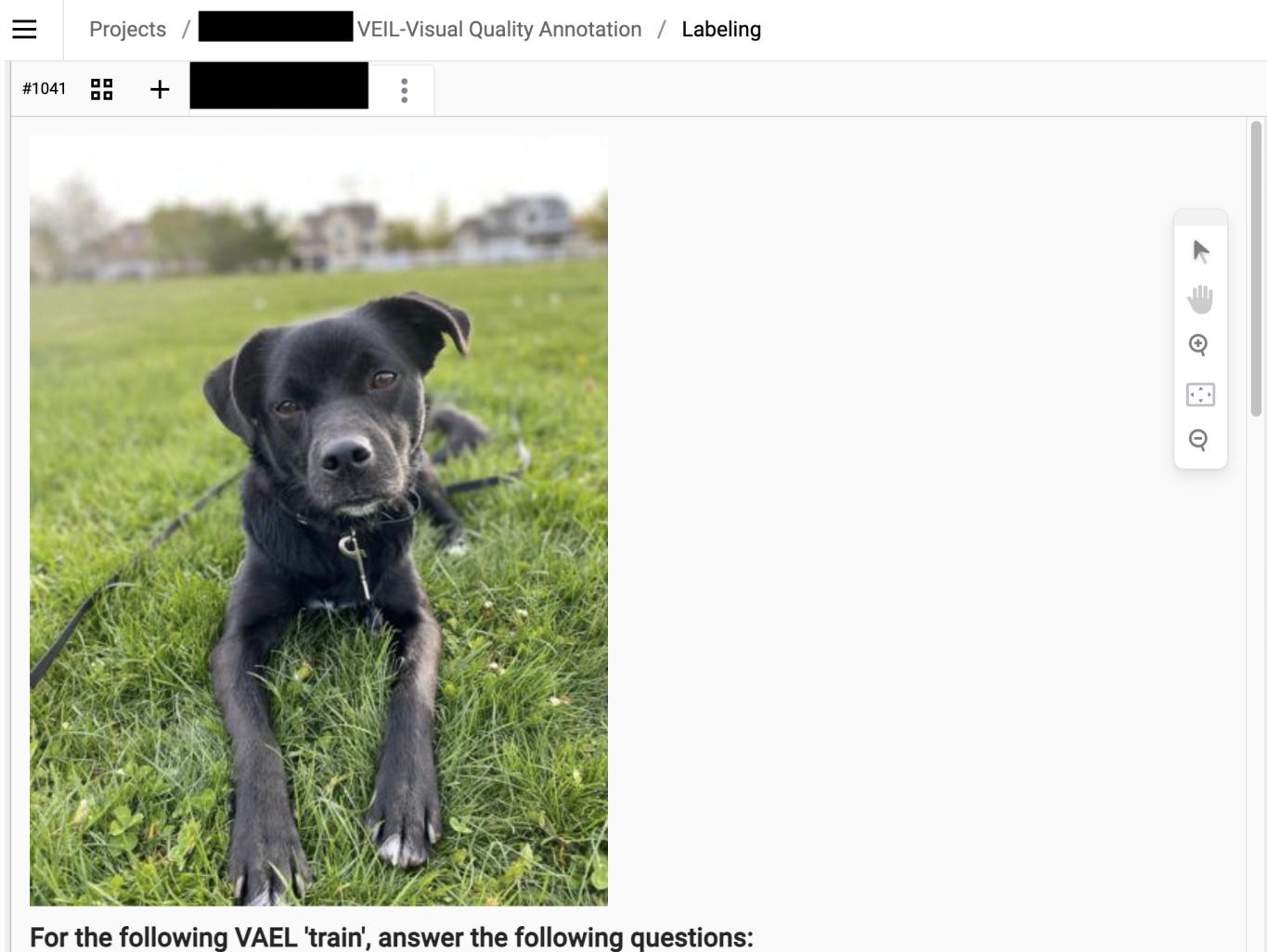
The screenshot shows the 'Cloud Storage' section of the Label Studio interface. On the left sidebar, 'Cloud Storage' is selected. The main area displays a single storage configuration:

Type	Local files
Path	/[REDACTED]/veil_annotation_project/[REDACTED]annotation
Status	Completed
Annotations	0 (0 total)
Last Sync	July 26, 2023 · 17:35:37

Buttons for 'Add Source Storage' and 'Add Target Storage' are visible. A 'Sync Storage' button is located at the bottom right of the storage card.

Should the images not load, then there might be an error with the relative paths in the export file. Check if you can construct the absolute file path by concatenating the output of `echo LABEL_STUDIO_LOCAL_FILES_DOCUMENT_ROOT` and `image_url` in the export json files (after `?d=`).

## Annotation Interface Example



For the following VAEI 'train', answer the following questions:

### Q1: Object Visibility (select which one is applicable)

VAEI: train

CAPTION: riley just came home after finishing his 2 week board and train!

- Option 1: Visible - 75% or more of object is visible<sup>[1]</sup>
- Option 2: Partially Visible - less than 75% of object is visible<sup>[2]</sup>
- Option 3: Completely Missing<sup>[3]</sup>
- Option 4: Not clear what fraction of the object is visible<sup>[4]</sup>

### Q2: If Option 3 [Object completely missing] was selected in Q1 (check all that apply)

VAEI: train

CAPTION: riley just came home after finishing his 2 week board and train!

- Image contains co-occurring scene/objects/entities with extracted label (boat and water, train and train tracks, bus and bus stop)<sup>[5]</sup>
- Semantically similar object is present in image instead (e.g. extracted label is truck but there is a car present in the image)<sup>[6]</sup>

### Q3: Object Appearance Quality - Defects (check all that apply)

VAEI: train

CAPTION: riley just came home after finishing his 2 week board and train!

- Occlusion (for example: inside a vehicle, or object blocked by another entity in an image)<sup>[7]</sup>
- Key Parts are missing from image (like head of an animal)<sup>[8]</sup>

Atypical appearance (clip art, watermark, knitted animal)<sup>[9]</sup>

Other<sup>[0]</sup>

If other is selected above, please specify here

#### **Q4: Check if the VAEI could be inferred to be visually absent due to proximity or membership in any of these linguistic elements (check one):**

VAEI: train

CAPTION: riley just came home after finishing his 2 week board and train!

Describing the Past<sup>[a]</sup>

Text discusses content/events related to but beyond the image<sup>[w]</sup>

Non-Literal (metaphor, simile, etc)<sup>[e]</sup>

Noun Modifier: Noun is modifying other noun (VAEI can be either noun)<sup>[t]</sup>

Prepositional Phrase (e.g. 'on a boat')<sup>[a]</sup>

Different Word Sense<sup>[s]</sup>

Named Entity<sup>[d]</sup>

Caption transcribes text in image<sup>[f]</sup>

Other<sup>[g]</sup>



Update

## Agreement

	Q1	Q2	Q3	Q4
RedCaps	0.75	0.14	0.29	0.56
SBUCaps	0.69	0.38	0.49	0.57
CC	0.83	0.33	0.52	0.61

Table: Annotator Agreement with Weighted Cohen's Kappa

## Agreement by Option Figures per Dataset

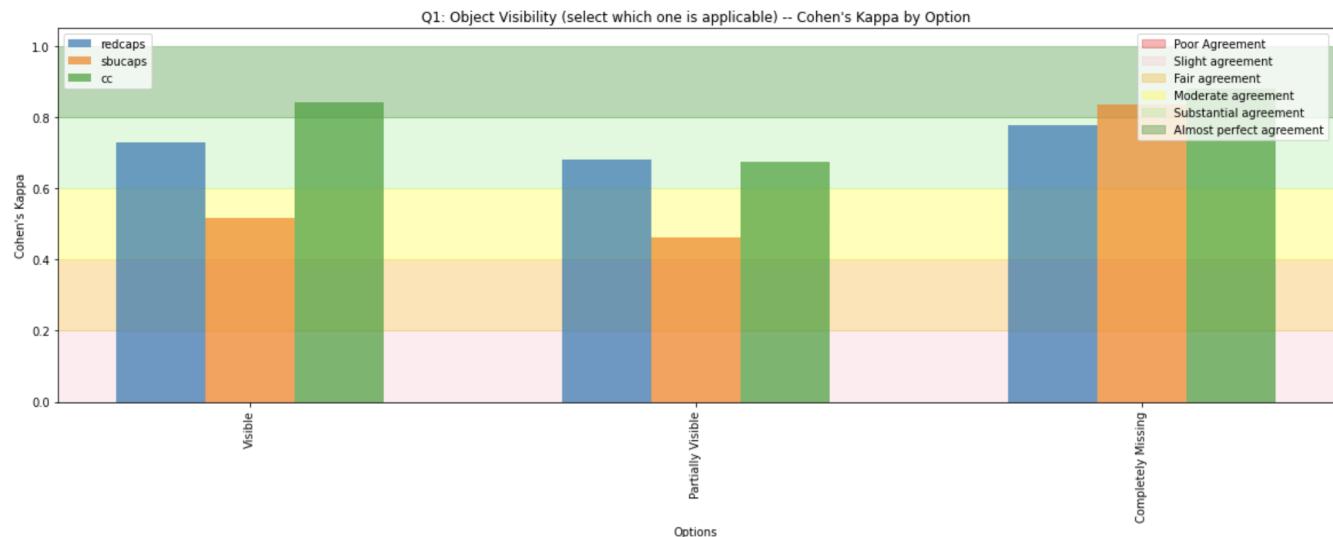


Figure: Annotator Agreement for Question 1

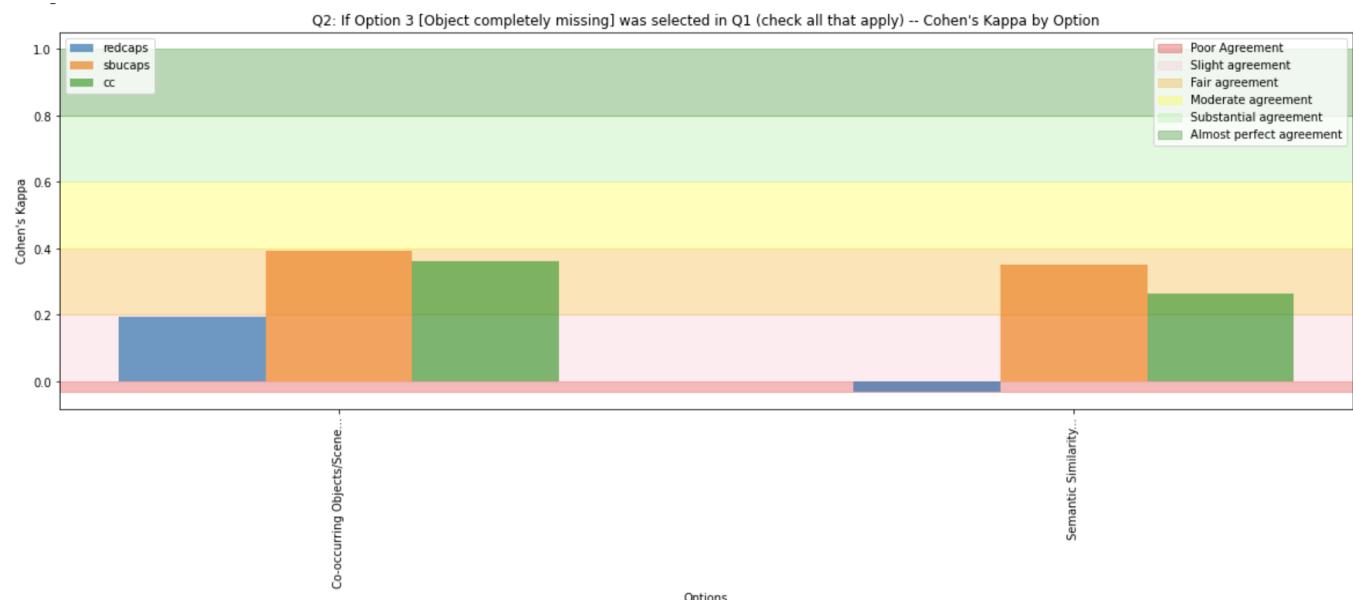


Figure: Annotator Agreement for Question 2

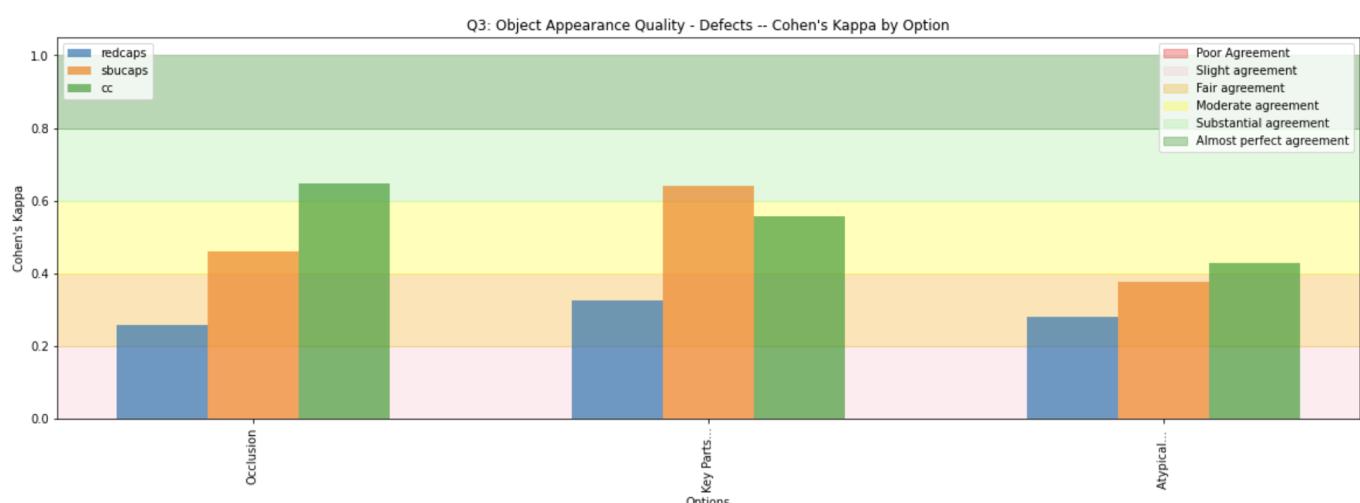


Figure: Annotator Agreement for Question 3

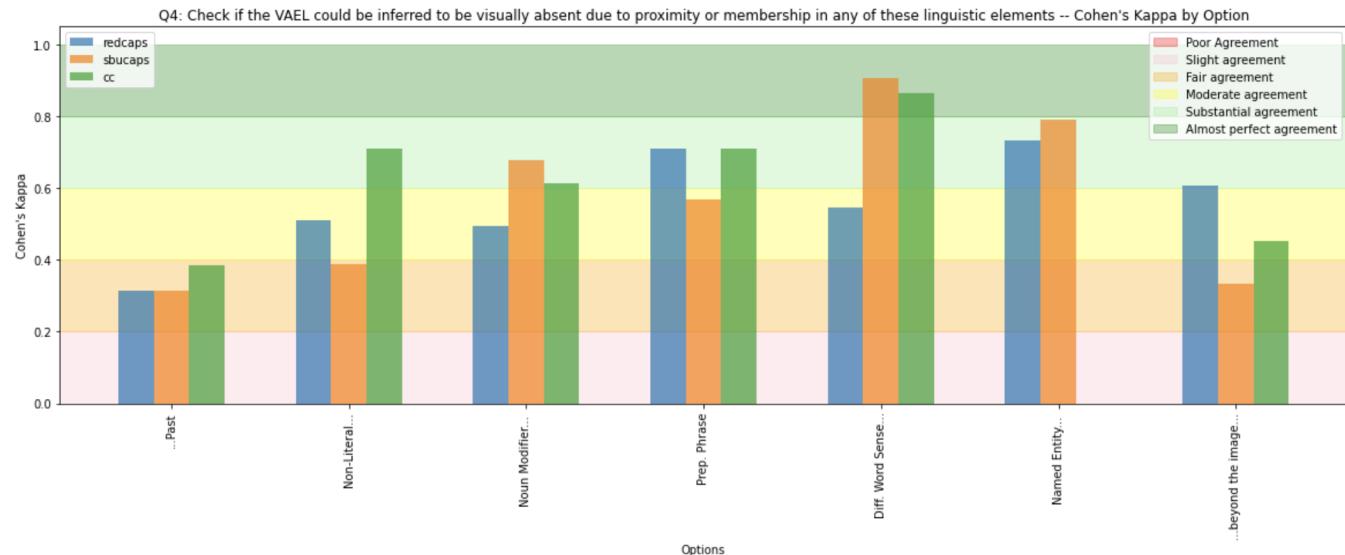


Figure: Annotator Agreement for Question 4