

Annotation Guidelines

English word/phrase 1

longer piece of English text 2

word/phrase does not exist 3

grammar mistake 4

spelling mistake 5

strange/wrong construction 6

strangely/wrongly used word/phrase 7

other linguistic remark [add to shared doc] 8

non-linguistic/meta remark [add to shared doc] 9

☐ Not sure^[o]

☐ Very minor issue/humans might write the same^[q]

☐ clearly from English^[w]

☐ could be from English^[e]

☐ no clear link to English^[t]

Na een stormachtige week waarin het bestuur van de prestigieuze Amerikaanse universiteit Harvard zich verdedigd heeft tegen ophef rond het behandelen van antisemitisme en beschuldigingen van plagiaat, heeft rector Lawrence S. Bacow aangekondigd dat hij niet meer terugkeert in zijn functie.

Bacows besluit komt na een interne onderzoekscommissie **die naar voren bleek** ^{strange/wrong construction} dat er fouten waren gemaakt bij het behandelen van meldingen over antisemitisme aan campus. De commissie concludeerde ook dat er sprake was geweest van "een gebrek aan leiderschap" op het hoogste niveau van de universiteit.

There are 9 different labels, which should be relatively self-explanatory. Remember that there are additional options once you have assigned a label.

- Generally it is important to try to be as objective as possible and not to impose your own style. Many sentences could be improved but are not actually wrong or even strangely written. Focus on real errors and things that you don't think native speakers would generally write.
- Select no more or no less of the text than is relevant for the label.
- (partially) overlapping annotations are possible

There are 2 options available for every annotation regardless of the label:

English word/phrase 1
longer piece of English text 2
word/phrase does not exist 3

grammar mistake 4
spelling mistake 5
strange/wrong construction 6

strangely/wrongly used word/phrase 7
other linguistic remark [add to shared doc] 8

non-linguistic/meta remark [add to shared doc] 9

☐ Not sure^[o]

☐ Very minor issue/humans might write the same^[q]

☐ clearly from English^[w]

☐ could be from English^[e]

☐ no clear link to English^[t]

Na een stormachtige week waarin het bestuur van de prestigieuze Amerikaanse universiteit Harvard zich verdedigd heeft tegen ophef rond het behandelen van antisemitisme en beschuldigingen van plagiaat, heeft rector Lawrence S. Bacow aangekondigd dat hij niet meer terugkeert in zijn functie.

Bacows besluit komt na een interne onderzoekscommissie **die naar voren bleek** ^{strange/wrong construction} dat er fouten waren gemaakt bij het behandelen van meldingen over antisemitisme aan campus. De commissie concludeerde ook dat er sprake was geweest van "een gebrek aan leiderschap" op het hoogste niveau van de universiteit.



If you have serious doubts about whether the annotation you gave is necessary or just personal preference, indicate the “not sure” option.



If you know the issue is very minor or something human native speakers would also write, then check the second option. This applies mostly to grammar and spelling and should not be used for the “English word/phrase” or “longer piece of English text” options.

English word/phrase

- not usually used in Dutch/French (e.g., self-explanatory, polar bear)
- sometimes used in Dutch/French (e.g., blueprint, air base)
- very commonly used in Dutch/French (e.g., blackmail, browser, database)

Additional info:

Anytime an English word or short phrase or multiword is used, annotate it with this label.

- Names that are not usually changed in Dutch do not count, unless there is a common Dutch alternative (e.g., do label “United States” because we have a common Dutch alternative “Verenigde Staten”, but don’t annotate “New York Academy of Art” because, though we can translate the name, there is no official Dutch version of the name.
- Words that have an English origin but have become so common in Dutch/French that we no longer really recognise them as English and/or the Dutch/French equivalent would sound very strange to use, do not have to be annotated with this label (e.g., computer, fitness, checken, etc.).
- The boundaries between the different options (not usually used, sometimes used, or very commonly used in Dutch/French) are of course quite subjective. Consider the context a bit: newspaper articles to be published on VRT or RTBF. If the English words are so commonly used that they could appear as such in these types of articles, then “very common” option is appropriate. If you think the English word is not known to the typical Belgian, then the “not usually used” option is probably appropriate.

longer piece of English text

- o entire text
- o part of text

Additional info:

Once more than a single word or short phrase is written in English, use this label instead of the previous one. Assuming that this is not that common in Dutch/French articles, you do not need to indicate how common it is, just whether it concerns the entire text or only part of it (though that should in principle also be clear from your selection of the text to be annotated).

- Select the entire text that is in English for this annotation and don’t make separate annotations unless there is Dutch/French text in between.
- If this occurs, you do not need to annotate this text any further. You can give it a cursory glance to check for very weird things that pop out and which can be mentioned in the shared spreadsheet, but otherwise you don’t have to spend more time and effort on these pieces of text.

word/phrase does not exist

- o literal ‘translation’ from English
- o probably influenced by English
- o no clear link to English

Additional info:

When words or phrases are used that simply don’t exist in French/Dutch and that are not simple spelling mistakes, use this label.

- Examples of literal ‘translations’: “een *zelfloze* daad” (~ selfless act), “op het blok” (~ on the block), “kubben” (~cubs).

grammar mistake

- clearly from English
- could be from English
- no clear link to English

Additional info:

Any grammatical error, e.g., “deze transport”, “het vis”, “gezonder levensstijl”. Can be harder to see a clear link to English.

spelling mistake

- clearly from English
- could be from English
- no clear link to English

Additional info:

Any spelling mistake, e.g., compound written in two words instead of one, wrongly place apostrophe, etc.

strange/wrong construction

- clearly from English
- could be from English
- no clear link to English

Additional info:

This category is for any turn of phrase that does not make sense. Often you can still tell what is meant, but it is just not how any native speaker would say it.

e.g., “kinderen van zo een jonge leeftijd”, “ter nagespeurd”, “kan een nieuw lag bereikt worden”

- This is one of the more difficult categories both in terms of determining whether something is really wrong, or just not how you would write it, and determining how much you should annotate (which words are part of the construction). Don't be afraid of starting a discussion with the other annotators when you're really unsure

Strangely/wrongly used word/phrase:

- clearly from English
- could be from English
- no clear link to English

Additional info:

This is very similar to the previous category, but at word/phrase level of course.

- Make sure not to use this instead of “word/phrase does not exist”
- Examples: “windjes” when talking about gusts of wind, “kennen” instead of “weten”.

other linguistic remark [add to shared doc]:

Additional info:

Anything that is clearly strange or wrong and has to do with the language, should be annotated with this label, and a remark should be made in the share spreadsheet to clarify. You can also use the comments in Label Studio, but please always make sure to use the spreadsheet as well.

Some examples based on things I have seen:

- text in German or Spanish
- writing that is clearly not the register of a newspaper article

non-linguistic/meta remark [add to shared doc]

Additional info:

Anything else that is strange but not related to the language can be indicated with this label.

Important: do not spend too much time and effort on this, as this is not the focus of this study. You are not supposed to be fact checking the texts. However, sometimes there are simply weird/funny things that feel strange not to annotate, or that can be interesting examples of characteristics of texts written by LLMs, so we did provide at least the option to annotate non-linguistic information

Some examples:

- Meta-information:
 - o The model writes that it does not know Dutch (in Dutch)
 - o A declaration is added that the article was written by a language model
 - o I wrote something between [] about manually stopping the model when it was stuck in a loop
- Clearly wrong information:
 - o Mentioning that penguins are mammals
- Other
 - o Article is just a repetition of the title