

## ANNOTATION GUIDELINES FOR THE BAR CHART DATASET OF PLOTS AND THEIR CROWDSOURCES SUMMARIES

Contributors: Rudy Khalil, Iza Škrjanec, Emeka Udeh, Vera Demberg

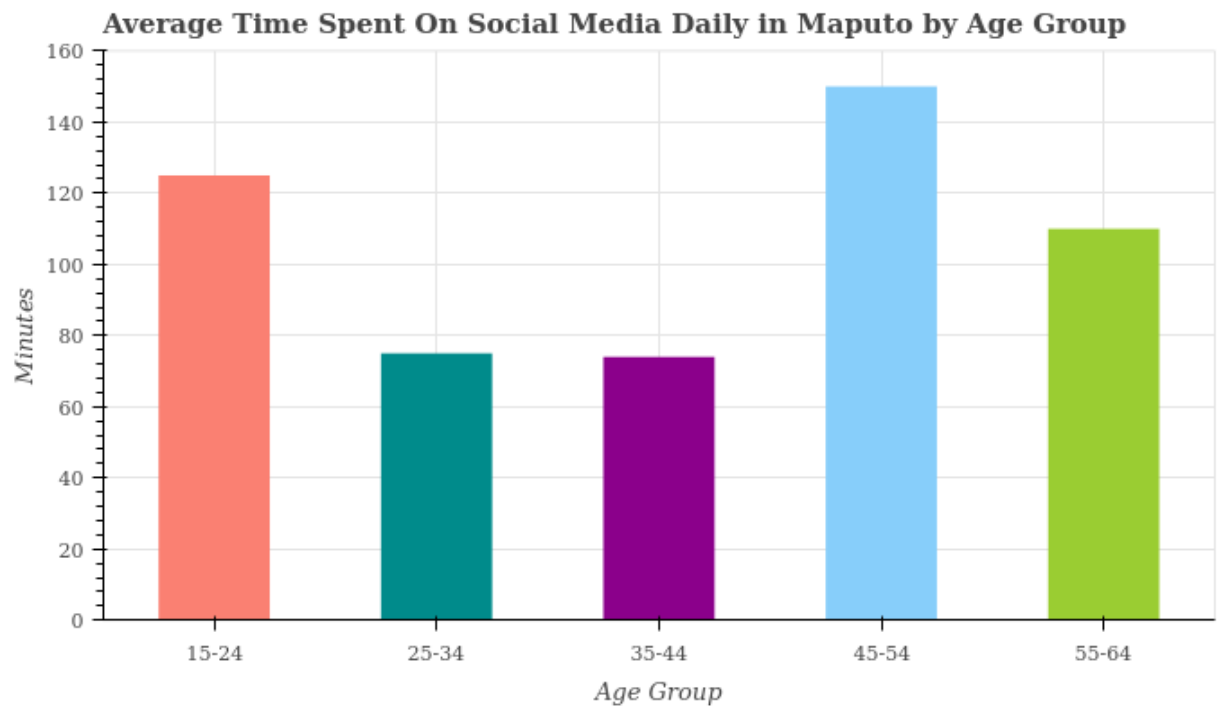
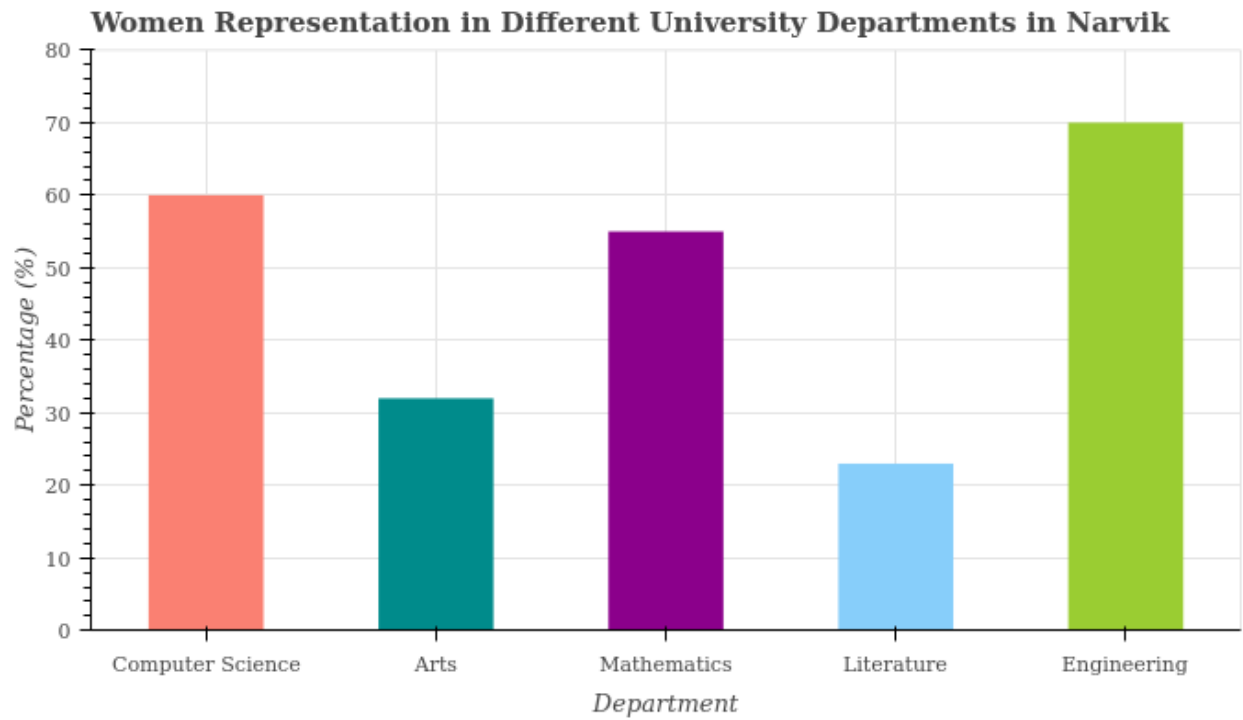
These guidelines present the set of labels used for annotating phenomena in summaries of bar charts. Previous versions of the guidelines can be found here: [version 1.0](#) and [version 1.1](#).

To collect data, we first generated bar chart plots and then collected their summaries via crowdsourcing. The current dataset has 47 plots, and for each of these, around 23 summaries.

When designing a tagging system, our main principle was to have a manageable set of labels which would help us 1) track discourse phenomena of chart summaries relevant to NLG, 2) delexicalize chart-specific tokens so as to make the input data more general, but also to relexicalize the output accordingly.

Each label roughly belongs to one of the following types of labels:

- 1) Chart-specific: names and heights of bars, other like the chart title and related vocabulary
- 2) Discourse: words/phrases typical of chart summaries, but not dependent on particular charts and their data
- 3) Interpretation: any kind of notions that are not grounded in the chart, be it interpretation of the data or non-relevant comment about the chart ("but the chart doesn't tell us which year the numbers were measured")



## 1) x\_axis

References to the name of the x axis is labeled as x\_axis. The reference may be verbatim as in the chart or taking a different surface form.

Examples for chart 1: department, departments, university departments, courses

Examples for chart 2: age group, age groups, group, groups, age range, user group

## 2) y\_axis

This label accounts for the quantity presented on the y axis. This may or may not overlap with the name of the y axis in the chart.

Chart 1: percentage\*, amount, share of women, women representation

Chart 2: time, time spent on social media, minutes\*

\*In case the name of the y axis and its unit overlap, pay attention to how a word is used in context. For example:

"The percentage [y\_axis] of women varies across different departments. In computer science it is 60 percent [y\_axis\_inferred\_label]."

## 3) x\_axis\_label \*\_value

References to bar names are annotated with these labels. Currently, we code each bar name in terms of its rank given the height. For charts with an ordinal X variable, we are considering coding in terms of rank given the x axis (order of appearance in the chart).

The label includes the rank information as follows:

- the highest bar: x\_axis\_label\_highest\_value
- second highest: x\_axis\_label\_Scnd\_highest\_value
- third highest: x\_axis\_label\_3rd\_highest\_value
- fourth highest: x\_axis\_label\_4th\_highest\_value
- ...
- least/lowest one: x\_axis\_label\_least\_value

Note that the label for the lowest bar doesn't include the rank number, but rather "least". If a chart has 4 bars of different heights, the subset of labels to be applied is {x\_axis\_label\_highest\_value, x\_axis\_label\_Scnd\_highest\_value, x\_axis\_label\_3rd\_highest\_value, x\_axis\_label\_least\_value}.

The reference to the bar may or may not match the bar name completely.

Chart 1, x\_axis\_label\_3rd\_highest\_value: mathematics, maths, math

Chart 2, x\_axis\_label\_Scnd\_highest\_value: 15-24, between 15 and 24, from 15 to 24

In case two or more bars are of exactly the same height, they should carry the same label.  
TODO check and give examples

#### 4) y\_axis\*\_value (\_val)

This label annotates bar heights. More specifically, the exact values of bar heights as they appear in the plotting data and in the chart itself.

The label names include the height rank:

- y\_axis\_highest\_value\_val
- y\_axis\_Scnd\_highest\_value
- y\_axis\_3rd\_highest\_value
- ...
- y\_axis\_least\_value\_val

Note that the label for highest and lowest heights have a suffix \_value\_val.

Chart 1:

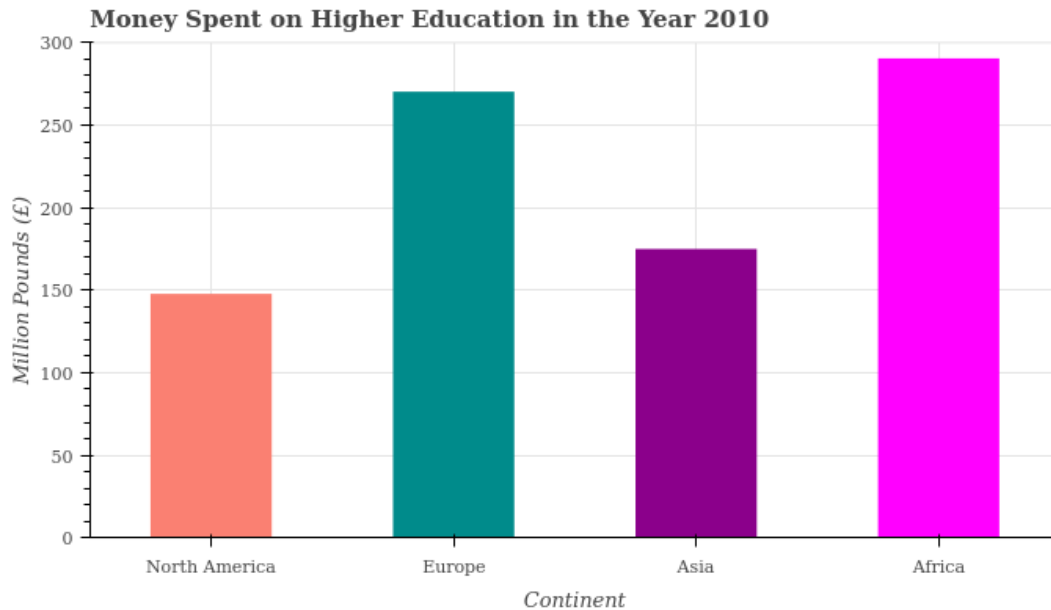
- y\_axis\_highest\_value\_val: 23
- y\_axis\_Scnd\_highest\_val: 32
- y\_axis\_3rd\_highest\_val: 55
- y\_axis\_4th\_highest\_val: 60
- y\_axis\_least\_value\_val: 70

For some charts, the y axis ticks mark lower integers, while the name of the y axis gives information about their magnitude (e.g. thousand or million). In that case, we still annotate the exact bar height, splitting the magnitude and labeling it separately.

For a description applying to the bottom chart, the labeling for this sentence is as follows:

*The amount of money spent on Higher Education in 2010 in Asia was £ 175,000,000.*

The  
amount of money   <y\_axis>  
spent  
on  
Higher Education in 2010   <topic\_related\_property>  
in  
Asia   <x\_axis\_label\_3rd\_highest\_value>  
was  
£   <y\_axis\_inferred\_label>  
**175**   <y\_axis\_3rd\_highest\_val>  
**,000,000**   <y\_magnitude>



TODO: if two bars are of the exact same height, do they share the label too?

## 5) y\_axis\_inferred\*\_value\_approx

Crowdsourcers often describe bars in terms of inexact heights. This could be because the top of the bar is not aligned with a tick mark stating the value. The bar heights are then inferred from the chart and given with an approximated value. Such instances are annotated with y\_axis\_inferred\*\_value\_approx.

The label names encode the height rank:

- y\_axis\_inferred\_highest\_value\_approx
- y\_axis\_inferred\_Scnd\_highest\_value\_approx
- ...
- y\_axis\_inferred\_least\_value\_approx

The

least <y\_axis\_least\_value>

amount of time <y\_axis>

was

spent

by

the

35-44 <x\_axis\_label\_least\_value>

age group <x\_axis>

totalling

**68 <y\_axis\_inferred\_least\_value\_approx> ## actual bar height: 74**

minutes <y\_axis\_inferred\_label>

For treating magnitudes, the same rule as for y\_axis\*\_value (\_val) applies.

## 6) y\_axis\_inferred\_label

The quantity on the y axis is measurable in particular units, which are often given as the whole or a part of the name of the y axis. The label used should be y\_axis\_inferred\_label.

The unit can be verbalized as a word ("pounds") or with a symbol ("£").

For chart 1, the label can be applied to e.g. "percents" or "%".

For chart 2, to "minutes", "min", "mins", "hours".

For chart 3, "pounds", "pound", "£".

Note that sometimes participants use a wrong currency, e.g. dollar instead of pound. We label such instances with y\_axis\_inferred\_label and correct the token, keeping both, the incorrect and corrected one.

## 7) y\_magnitude

To keep the chart clear and tick marks on the y axis visible, the tick marks are often smaller numbers, while the name of the y axis gives information about their magnitude.

See chart 3 as an example. The highest tick mark goes to 300, the magnitude is million.

Magnitudes can be verbalized in different forms:

- thousand: "thousand", "thousands", "k", "000", ",000"
- million: "million", "millions", "m", "M", "000000", ",000,000"

If the bar height appears together with the magnitude, we split them and label accordingly.

See below:

*It actually decreased to nearly £ 19,500 by the year 2010.*

it

actually

decreased <slope\_down>

to

nearly <y\_axis\_approx>

£ <y\_axis\_inferred\_label>

19.5 <y\_axis\_inferred\_3rd\_highest\_value\_approx>

**000** <y\_magnitude>

by

the

year <x\_axis>

2010 <x\_axis\_label\_3rd\_highest\_value>

*last was North America who only spent nearly £ 150 m on higher education.*

last <order\_last>

was

North America <x\_axis\_label\_least\_value>

who

only <y\_axis\_least\_value>

spent  
 nearly <y\_axis\_approx>  
 £ <y\_axis\_inferred\_label>  
 150 <y\_axis\_inferred\_least\_value\_approx>  
**m** <y\_magnitude>  
 on  
 higher education <topic\_related\_property>

## 8) y\_axis\_inferred\_value\_add\_v1=\*\_v2=\*

When the relative relation between bar heights is calculated via addition/subtraction, we label the height value as y\_axis\_inferred\_value\_add.

The label for additive relation assumes two bar heights are involved, v1 and v2. The relative bar height d is calculated, such that  $v1 - v2 = d$ . The given relative height does not have to be exactly equal to d; it may be an approximation.

The attributes v1 and v2 get as values the names of bars, which are coded in terms of their height rank. For example, y\_axis\_inferred\_value\_add\_v1=Scnd\_v2=4th annotates the relative height that we get from height\_Scnd minus height\_4th.

*Insurance has the highest representation at 65% whereas law firms has the lowest at 35% - a 30% difference.*

Insurance <x\_axis\_label\_highest\_value>  
 has  
 the  
 highest <y\_axis\_highest\_value>  
 representation <y\_axis>  
 at  
 65 <y\_axis\_highest\_value\_val>  
 % <y\_axis\_inferred\_label>  
 whereas <y\_x\_comparison>  
 law firms <x\_axis\_label\_least\_value>  
 has  
 the  
 lowest <y\_axis\_least\_value>  
 at  
 35 <y\_axis\_least\_value\_val>  
 % <y\_axis\_inferred\_label>  
 -  
 a  
**30** <y\_axis\_inferred\_value\_add\_v1=highest\_v2=least>  
 % <y\_axis\_inferred\_label>  
 difference

*Those in the 15-24 age bracket spent the most time on social media, clocking over 100 minutes more than those in the 55-64 age bracket.*

Those  
in  
the  
15-24 <x\_axis\_label\_highest\_value>  
age bracket <x\_axis>  
spent  
the  
most <y\_axis\_highest\_value>  
time on social media <y\_axis>  
,  
clocking  
over <y\_axis\_approx>  
**100** <y\_axis\_inferred\_value\_add\_v1=highest\_v2=least>  
minutes <y\_axis\_inferred\_label>  
more <y\_axis\_trend\_up>  
than  
those  
in  
the  
55-64 <x\_axis\_label\_least\_value>  
age bracket <x\_axis>  
.

*the following group 25-34 spend about 160 minutes, 20 less that the previous group*  
the  
following <order\_Scnd>  
group 25-34 <x\_axis\_label\_Scnd\_highest\_value>  
spend  
about <y\_axis\_approx>  
160 <y\_axis\_inferred\_Scnd\_highest\_value\_approx>  
minutes <y\_axis\_inferred\_label>  
,  
**20** <y\_axis\_inferred\_value\_add\_v1=Scnd\_v2=highest>  
less <y\_axis\_trend\_down>  
that  
the  
previous group <x\_axis\_label\_highest\_value>

## 9) y\_axis\_inferred\_value\_mul\_v1=\*\_v2=\*

When the relative relation between bar heights is calculated via multiplication/division, we label the height value as y\_axis\_inferred\_value\_mul.

Similarly as for \_add\_, we assume two bars are involved in the relation. The relative height k is calculated as follows:  $v1 / v2 = k$ . The given relative height does not have to be exactly equal to k; it may be an approximation.



Note that k can be expressed as a number (2), a word (twice) or a combination (2 times).

*STEM subjects are half as popular as Philosophy with 15% opting for this.*

STEM subjects <x\_axis\_label\_least\_value>  
are  
**half** y\_axis\_inferred\_value\_mul\_v1=least\_v2=Scnd>  
as  
popular  
as  
Philosophy <x\_axis\_label\_Scnd\_highest\_value>  
with  
15 <y\_axis\_least\_value\_val>  
% <y\_axis\_inferred\_label>  
opting  
for  
this

*In law firms this is 3 times higher than in tech where only 20% are women.*

In  
law firms <x\_axis\_label\_highest\_value>  
This  
is  
**3 times** <y\_axis\_inferred\_value\_mul\_v1=highest\_v2=least>  
higher <y\_axis\_trend\_up>  
than  
in  
tech <x\_axis\_label\_least\_value>  
where  
only  
20 <y\_axis\_least\_value\_val>  
% <y\_axis>  
are  
woman <topic\_related\_property>

## 10) slope\_x\_value

When the variable on the x axis is interval/ratio or even ordinal, the slope of the regression line can be calculated and described in the chart summary.

Typically, the label slope\_x\_value tells us the step size on the x axis in the slope expression. The slope value usually applies to the entire data, not a subset. The given step size can be exact or an approximation.

Take for example the following data showing the number of fatal injuries in the span of 5 years:

- x: '2012', '2013', '2014', '2015', '2016'
- y: 30, 25, 16, 15, 12

- We can see that the slope is negative; the regression line is falling; the number of injuries drops year by year
- In the example sentence below, the slope is 5 injuries per 1 year. Note that "year" here means "one year", so it will be annotated this way.

*The trend decreases from 30 in 2012 to 13 in 2016, with an average decrease of 5 per year.*

The  
trend <y\_axis\_trend>  
decreases <slope\_down>  
from  
30 <y\_axis\_highest\_value\_val>  
in  
2012 <x\_axis\_label\_least\_value>  
to  
13 <y\_axis\_inferred\_least\_value\_approx>  
in  
2016 <x\_axis\_label\_least\_value>  
,  
with  
an  
average <y\_axis\_trens>  
decrease <slope\_down>  
of  
**5** <slope\_y\_value>  
per  
**year** <slope\_x\_value>

## 11) slope\_y\_value

See slope\_x\_value. Similarly, we annotate the step size on the y axis. The size can be exact or an approximation

In the example above, the step size is -5: 5 death less per year.

*The trend decreases from 30 in 2012 to 13 in 2016, with an average decrease of 5 per year.*

## 12) x\_axis\_range\_start and x\_axis\_range\_end

The markers of spans between bars are labeled with the range start and range end label.

This applies especially to charts with a temporal variable on the x-axis.

The typical pairs of start and end markers are: between-and, from-to, from--

The  
significant  
results  
show

that  
the  
number of fatal injuries <y\_axis>  
has  
been  
declining <slope\_down>  
**from** <x\_axis\_range\_start>  
year <x\_axis>  
2012 <x\_axis\_label\_highest\_value>  
up  
**to** <x\_axis\_range\_end>  
year <x\_axis>  
2016 <x\_axis\_label\_least\_value>

the  
number of fatal injuries at the Pula Steel factory <topic>  
**from** <x\_axis\_range\_start>  
2012 <x\_axis\_label\_highest\_value>  
- <x\_axis\_range\_end>  
2016 <x\_axis\_label\_least\_value>

### 13) x\_interval

The label x\_interval can apply to charts with an ordinal/interval/ratio variable on the x-axis. It is currently used in 3 different contexts:

- a) Referring to the total interval on the x\_axis (between the first and the last bar)
- b) Referring to the interval between two bars
- c) Referring to the interval within a bar (for bars representative of time interval)

Note that it does not annotate slope values.

Example for a)

*This chart represents the representation of minorities in the Parliament of Libya over two decades from 1990 to 2019.*

This  
Chart  
represents  
the  
representation of minorities in the Parliament of Libya <topic>  
over  
**two decades** <x\_interval>  
from <x\_axis\_range\_start>  
1990 <x\_axis\_label\_highest\_value>

to <x\_axis\_range\_end>  
2019 <x\_axis\_label\_5th\_highest\_value>

Example for b)

*The graph shows that the median salaries of women have risen each half decade.*

The  
graph  
shows  
that  
the  
median salaries of women <topic>  
have  
risen <slope\_up>  
each  
**half decade** <x\_interval>

Example for c)

*This graph shows the average time spent on social media by age group in Maputo.  
Five 10 year spans from 15 to 64 are given.*

This  
graph  
shows  
the  
average time spent on social media by age group in Maputo <topic>  
.  
Five <x\_axis\_labels\_count>  
**10 year** <x\_interval>  
spans <x\_axis>  
from <x\_axis\_range\_start>  
15 <x\_axis\_label\_highest\_value>  
to <x\_axis\_range\_end>  
64 <x\_axis\_label\_least\_value>  
are  
given

## 14) x\_axis\_labels\_count

Explicit references to the number of bars in the chart are labeled as x\_axis\_labels\_count. This applies only to the exact number of bars; whereby the numeral can be written as an integer ("5") or spelled out ("five").

The  
graph  
details  
the

gender pay gap <y\_axis>  
in  
**three** <x\_axis\_labels\_count>  
European <t=0\_a=x>  
countries <x\_axis>

## 15) topic

The topic label largely fits the title of the chart. The title often includes the names of the x- and y-axis, and their units, but this varies across charts. This typically appears at the beginning of a description as an introductory sentence.

It is important to note that the topic label encompasses only what appears in the title and the synonyms.

Chart title: *Women representation in different university department in Narvik*

The  
chart  
shows  
the  
percentage <y\_axis>  
of  
**women in different university departments in Narvik** <topic>

Chart title: *What causes obesity in Kiribati*

This  
chart  
looked  
at  
**causes of Obesity in Kiribati** <topic>

Chart title: *Women representation in different sectors in Benoni*

The  
following  
graph  
shows  
the  
**female representation in different work sectors in Benoni** <topic>

Note that the title is often broken down into parts in the description. The parts are then labeled accordingly as topic or with other fitting labels.

Chart title: *Average time spent on social media daily in Maputo by age group*

**Daily social media use by time in Maputo** <topic>  
is  
shown

here  
by  
age group <x\_axis>

## 16) group\_\* [in progress]

This label can account for cases when participants join bars and talk about them as a group. The label should include the names of bars which are inside the group. The names have been coded given their height so far.

For example, group\_Scnd\_3rd\_4th is a group consisting of the second, third and the fourth highest bar.

So far, the second batch of data has shown examples of:

- a) **conceptual grouping**: the bars share some hypernym-like concept. For example, "lifestyle choices" (referring to the bars of *Fast Food* and *Lack of Exercise*, as opposed to *Genetics* as causes of obesity)
- b) **grouping by similar height**: bars have a similar height, so their descriptions share a predicate.
- c) **grouping by different heights**: based on our observations, this applies to bar charts with an ordinal X variable. So far, we've seen the following case: there was a similarly large height difference between 4 bars (2 bars at the beginning of the x axis, 2 bars at the end). This is the single case so far, so we might not stick to it.
- d) **\*\* one versus rest**: highest/lowest bar versus rest: this grouping might appear when the isolated bar has a very different height than the rest of the bars.
- e) **Bars of exact same height**:

Example for b):

*The lowest users by time are in the 25-34 and 35-44 age groups, which use social media for around 75 minutes each day on average.*

The  
lowest <y\_axis\_least\_value>  
users <x\_axis>  
by  
time <y\_axis>  
are  
in  
the  
**25-34 and 35-44** <group\_4th\_least>  
age groups <x\_axis>  
,  
which  
use  
social media <topic\_related\_object>  
for  
around <y\_axis\_approx>

75 <group\_y\_4th\_least>  
minutes <y\_axis\_inferred\_label>  
each day ##  
on average

Another example for b)

*It indicates a downward trend in fatalities between 2013 and 2015.*

It  
indicates  
a  
downward <slope\_down>  
trend  
in  
fatalities <y\_axis>  
between 2013 and 2015<group\_highest\_3rd\_least>

Example for d) based on [this chart](#):

*We can see that a high percentage choose STEM whilst a much lower percentage choose literature, philosophy and medicine.*

We  
can  
see  
that  
a  
high <y\_axis\_highest\_value>  
percentage <y\_axis>  
choose  
STEM <x\_axis\_label\_highest\_value>  
whilst <y\_x\_comparison>  
a  
much <y\_axis\_trend>  
lower <y\_axis\_least\_value>  
percentage <y\_axis>  
choose  
literature, philosophy or medicine <group\_Scnd\_3rd\_least>

Example for e)

*with literature and medicine both at 20%*

with  
**literature and medicine** <group\_Scnd\_3rd>  
both  
at  
**20** <group\_y\_Scnd\_3rd>  
% <y\_axis>

## 17) group\_y\_\*[in progress]

When a subset of bars is treated as a group, their heights are usually described in a single predicate as well. This group height can be the mean of heights or even a sum.

The label name is similar to the group\_\* label with the only difference that it indicated the y: group\_y\_\*

We are not sure yet if the label name should also include information on how the group height is calculated: mean, sum or something else.

In the example for the group\_\* label, the group height is the mean.

*The lowest users by time are in the 25-34 and 35-44 age groups, which use social media for around 75 minutes each day on average.*

## 18) y\_mean

The label marks references to the mean height of all the bars in the plot.

*Fatal injuries at the Pula Steel Factory have varied from 2012 to 2016, averaging above 25.*

Fatal injuries <y\_axis>  
at  
the  
Pula Steel Factory <t=1\_a=b>  
have  
varied <y\_axis\_trend>  
from <x\_axis\_range\_start>  
2012 <x\_axis\_label\_Scnd\_highest\_value>  
to <x\_axis\_range\_end>  
2016 <x\_axis\_label\_4th\_highest\_value>  
,  
averaging  
above <y\_axis\_approx>  
**25** <y\_mean>

## 19) y\_axis\_highest\_value

The vocabulary (words or phrases) describing the highest bar without providing the bar height are labeled as y\_axis\_highest value. This includes expressions of superlatives and majority.

Some examples include: "highest", "largest", "biggest", "most popular/common/frequent".



Note that the phrases might get between labels. In general, determiners are not included into the label span. For example, "the most common cause of obesity" is a noun phrase, which will be labeled as follows:

- the (no label)
- most common <y\_axis\_highest\_value>
- cause of obesity <x\_axis>

Note that this label annotates expressions that pre-modify bar names, the name of the y axis, or some topic-related entity.

More examples

*most young people spend their evenings reading a book*  
most <y\_axis\_highest\_value>  
young people <topic\_related\_property>  
spend  
their evenings <topic\_related\_property>  
reading a book <x\_axis\_label\_highest\_value>

*the most popular choice is reading a book*  
The  
most popular <y\_axis\_highest\_value>  
choice <x\_axis>  
is  
reading a book <x\_axis\_label\_highest\_value>

*Africa has the highest amount at £ 290 million*  
Africa <x\_axis\_label\_highest\_value>  
has  
the  
highest <y\_axis\_highest\_value>  
amount <y\_axis>  
at  
£ <y\_axis\_inferred\_label>  
290 <y\_axis\_inferred\_highest\_value\_approx>  
million <y\_magnitude>

## 20) y\_axis\_least\_value

The lowest bar or its height reference may be modified by an expression that is labeled with y\_axis\_least\_value. Such expressions include "lowest", "smallest", "least common/popular".

See y\_axis\_highest\_value for general principles for the span of the label.  
Some examples:

*in 2015 it was the lowest at £ 16 k*  
in

2015 <x\_axis\_label\_least\_value>  
it  
was  
the  
lowest <y\_axis\_least\_value>  
at  
£ <y\_axis\_inferred\_label>  
16 <y\_axis\_least\_value\_val>  
k <y\_magnitude>

*Literature has the fewest women*  
Literature <x\_axis\_label\_least\_value>  
has  
the  
fewest <y\_axis\_least\_value>  
women <topic\_related\_object>

As it has been observed in the data, the bar height of the lowest bar is often pre-modified by "only" or similar expressions. We label them as y\_axis\_least\_value as well.

*the group aged 35-44 with only 65 minutes a day*  
the  
group <x\_axis>  
aged  
35-44 <x\_axis\_label\_least\_value>  
with  
only <y\_axis\_least\_value>  
65 <y\_axis\_inferred\_least\_value\_approx>  
minutes <y\_axis\_inferred\_label>  
a day <topic\_related\_property>

## 21) y\_axis\_approx

Hedge expressions (e.g. "about", "around", "approximately", "more than", "less than", "over", "under", "just over") explicitly signal approximations. When these words play the role of a hedge, we label them as y\_axis\_approx.

The affiliated y values may or may not be inexact or rounded.

Some examples

*Roughly 31% of students chose to study medicine*  
Roughly <y\_axis\_approx>  
31 <y\_axis\_highest\_value\_val>  
% <y\_axis\_inferred\_label>  
of  
students <topic\_related\_property>

chose  
to <x\_axis\_range\_end>  
study  
medicine <x\_axis\_label\_highest\_value>

*and fast food the lowest, at just under 25%*  
and  
fast food <x\_axis\_label\_least\_value>  
the  
lowest <y\_axis\_least\_value>  
,  
at  
just under <y\_axis\_approx>  
25 <y\_axis\_inferred\_least\_value\_approx>  
% <y\_axis\_inferred\_label>

*Between 2000 and 2004 minorities had more than 7% representation.*  
Between 2000 and 2004 <x\_axis\_label\_4th\_highest\_value>  
minorities <topic\_related\_property>  
had  
more than <y\_axis\_approx>  
7 <y\_axis\_inferred\_4th\_highest\_value\_approx>  
% <y\_axis\_inferred\_label>  
representation <y\_axis>

## 22) y\_axis\_trend, y\_axis\_trend\_up, y\_axis\_trend\_down

When two or more bars are compared in a more vague manner, i.e. with adjectives and adverbs instead of providing their heights, one of the three labels is typically used: y\_axis\_trend, y\_axis\_trend\_up, y\_axis\_trend\_down.

Note that the label name is misleading: the label does not signify trends in terms of a regression line, but rather comparisons.

If a comparison has no direction, but just expresses a big, small, (in)significant difference without stating which bar is higher, the label used should be **y\_axis\_trend**.

*In that year the gender pay gap is similar in Germany and the UK at just over 20%.*  
In  
that  
year  
the  
gender pay gap <y\_axis>  
is  
**similar** <y\_axis\_trend>  
in

Germany <x\_axis\_label\_highest\_value> ## group\_highest\_Scnd  
 and  
 the  
 UK <x\_axis\_label\_Scnd\_highest\_value>  
 at  
 just over <y\_axis\_approx>  
 20 <y\_axis\_inferred\_Scnd\_highest\_value\_approx>  
 % <y\_axis\_inferred\_label>

*Mathematics follows closely with 56%*  
 Mathematics <x\_axis\_label\_3rd\_highest\_value>  
 follows <order\_3rd>  
**closely** <y\_axis\_trend>  
 with  
 56 <y\_axis\_3rd\_highest\_val>  
 % <y\_axis\_inferred\_label>

*This increases largely with those between age 45-54 spending the most time on social media*  
 This  
 increases <slope\_up>  
**largely** <y\_axis\_trend>  
 with  
 those  
 between age 45-54 <x\_axis\_label\_highest\_value>  
 spending  
 the  
 most <y\_axis\_highest\_value>  
 time <y\_axis>  
 on  
 social media <topic\_related\_object>

The label y\_axis\_trend\_(up|down) signifies base or comparative forms of adjectives/adverbs denoting size (in some way). Greatness is labeled with **y\_axis\_trend\_up**, e.g. "higher", "bigger", "larger", "more popular". Smallness (e.g. "smaller", "lower", "less common") with **y\_axis\_trend\_down**. Both of these labels can be modified by y\_axis\_trend.

Note that in case of phrases like "more than" and "less than", only the adverb is labeled. We leave "than" unlabeled; in the opposite case a new/different label would have to be used for "than" when the phrase is split. Compare "*Asia spends more than North America*" and "*Asia spends more money than North America*".

*Asia was the bigger of the two with around 170 millions.*  
 Asia <x\_axis\_label\_3rd\_highest\_value>  
 was  
 the

**bigger** <y\_axis\_trend\_up>  
 of  
 the  
 two  
 with  
 around <y\_axis\_approx>  
 170 <y\_axis\_inferred\_3rd\_highest\_value\_approx>  
 millions <y\_magnitude>

*The chart shows that there is more money spent on Higher Education in Asia than in North America*

The  
 chart  
 shows  
 that  
 there  
 is  
**more** <y\_axis\_trend\_up>  
 money <y\_axis>  
 spent  
 on  
 Higher Education <topic\_related\_object>  
 in  
 Asia <x\_axis\_label\_highest\_value>  
 than  
 North America <x\_axis\_label\_Scnd\_highest\_value>

*The amount in Europe was a little lower at £ 270,000,000*

The  
 amount <y\_axis>  
 in  
 Europe <x\_axis\_label\_Scnd\_highest\_value>  
 was  
 a  
 little <y\_axis\_trend>  
**lower** <y\_axis\_trend\_down>  
 at  
 £ <y\_axis\_inferred\_label>  
 270 <y\_axis\_inferred\_Scnd\_highest\_value\_approx>  
 ,000,000 <y\_magnitude>

Note that the same expression in different contexts carries a different role and thus label.  
 For example, "more than":

- y\_axis\_trend\_up: Asia has spent more than Europe.
- y\_axis\_approx: Africa has spent more than £ 60 millions on higher education.

## 23) slope\_up, slope\_down and slope\_mix

The labels slope\_[up|down] annotate lexical expressions of rising or falling slope in descriptions of ordinal/interval/ratio bars.

*2010-2014 saw a 6.2% minority and this increased to 7% between 2015-2019.*

2010-2014 <x\_axis\_label\_least\_value>  
saw  
a  
6.2 <y\_axis\_inferred\_least\_value\_approx>  
% <y\_axis\_inferred\_label>  
minority <y\_axis>  
and  
this  
**increased** <slope\_up>  
to  
7 <y\_axis\_inferred\_5th\_highest\_value\_approx>  
% <y\_axis\_inferred\_label>  
between 2015-2019 <x\_axis\_label\_5th\_highest\_value>

*On Monday the stock closed at £50, then rose to £63 on Tuesday.*

On  
Monday <x\_axis\_label\_3rd\_highest\_value>  
the  
stock <t=1\_a=y>  
closed <t=1\_a=y>  
at  
£ <y\_axis\_inferred\_label>  
50 <y\_axis\_inferred\_3rd\_highest\_value\_approx>  
,  
then <order\_Scnd>  
**rose** <slope\_up>  
to  
£ <y\_axis\_inferred\_label>  
63 <y\_axis\_inferred\_Scnd\_highest\_value\_approx>  
on  
Tuesday <x\_axis\_label\_Scnd\_highest\_value>

*It shows that there was a steady decline in injuries from 30 people being fatally injured in 2012 to 12 people being injured in 2016.*

It  
shows

that  
there  
was  
a  
steady <y\_axis\_trend>  
**decline** <slope\_down>  
in  
injuries <y\_axis>  
from  
30 <y\_axis\_highest\_value\_val>  
people being fatally injured <y\_axis\_inferred\_label>  
in  
2012 <x\_axis\_label\_highest\_value>  
to  
12 <y\_axis\_least\_value\_val>  
people being injured <y\_axis\_inferred\_label>  
in  
2016 <x\_axis\_label\_least\_value>

References to a dynamic trend of the bar heights are labeled as slope\_mix.

*The salary of women in Najaf fluctuated between £ 16,000 to £ 26,000*

The  
salary of women in Najaf <y\_axis>  
**fluctuated** <slope\_mix>  
between  
£ <y\_axis\_inferred\_label>  
16 <y\_axis\_least\_value\_val>  
,000 <y\_magnitude>  
to  
£ <y\_axis\_inferred\_label>  
26 <y\_axis\_inferred\_highest\_value\_approx>  
,000 <y\_magnitude>

## 24) x\_axis\_labels\_rest

When bars are split and referred to as one-versus-rest or some-versus-rest, the label x\_axis\_labels\_rest is used for the reference to the rest bars. Note that while this is a case of grouping, it is not necessarily motivated by conceptual similarity or height.

*Female representation in the Insurance sector in Benoni is higher than in Law Firms, Tech and Financial Groups.*

Female representation <y\_axis>  
in

the  
 Insurance sector <x\_axis\_label\_highest\_value>  
 in  
 Benoni <t=1\_a=b>  
 is  
 higher <y\_axis\_trend\_up>  
 than  
 in  
 Law firms , Tech and Financial Groups <x\_axis\_labels\_rest>

## 25) x\_axis\_trend

In charts with an ordinal/interval/ratio variable on the x-axis, descriptions of rising (or in rare cases, falling; x\_axis\_trend\_down) of this variable are annotated with x\_axis\_trend. The rising/falling refers to the bars' value on the x-axis, not on the y-axis.

[week days on x-axis] *Overall we can see a gradual decrease in closing stock across the week*

Overall  
 we  
 can  
 see  
 a  
 gradual <y\_axis\_trend>  
 decrease <slope\_down>  
 in  
 closing stock <y\_axis>  
**across the week <x\_axis\_trend>**

[age groups, younger to older on x-axis] *The older the group the less time they spend daily on social media*

**The older the group <x\_axis\_trend>**  
 he  
 less <slope\_down>  
 time <y\_axis>  
 they  
 spend  
 daily <t=1\_a=y>  
 on  
 social media <t=1\_a=y>

Note that, these expressions are realized as steps ("year on year") or frequency ("each year").

[years on x-axis] *The fatality rate has fallen each year.*  
 The



fatality rate <y\_axis>  
has  
fallen <slope\_down>  
**each year** <x\_axis\_trend>

Describing the bars in terms of right-to-left orientation on the x-axis should be labeled as x\_axis\_trend\_down.

[age groups in ascending order on x-axis] *The chart shows that the younger you are the more time you spend on social media.*

The  
chart  
shows  
that  
the  
**younger you are** <x\_axis\_trend\_down>  
the  
more <slope\_up>  
time <y\_axis>  
spent  
on  
social media <t=1\_a=y>

## 26) order\_\*

The label order\_\* applies to words/phrases that mark the order of bars, in terms of how the bars appear on the chart or in the description.

The order is explicit in the annotated word/token ("second", "thirdly") or inferred from context ("the highest is Spain, followed by Germany").

The label set includes order\_Sncd, order\_3rd, order\_4th, ... order\_last.

*Insurance is the highest at 64%, tech next at 61%, then financial groups at 50%, then law at 35%*

Insurance <x\_axis\_label\_highest\_value>  
is  
the  
highest <y\_axis\_highest\_value>  
at  
64 <y\_axis\_inferred\_highest\_value\_approx>  
% <y\_axis\_inferred\_label>  
,  
tech <x\_axis\_label\_Scnd\_highest\_value>

**next** <order\_Scnd>  
 at  
 61 <y\_axis\_inferred\_Scnd\_highest\_value\_approx>  
 % <y\_axis\_inferred\_label>  
 ,  
**then** <order\_3rd>  
 financial groups <x\_axis\_label\_3rd\_highest\_value>  
 at  
 50 <y\_axis\_3rd\_highest\_val>  
 % <y\_axis\_inferred\_label>  
 ,  
**then** <order\_last>  
 law <x\_axis\_label\_least\_value>  
 at  
 35 <y\_axis\_least\_value\_val>  
 % <y\_axis\_inferred\_label>

*The third is Asia which drops significantly compared to Europe*

The  
**third** <order\_3rd>  
 is  
 Asia <x\_axis\_label\_3rd\_highest\_value>  
 which  
 drops <y\_axis\_trend\_down>  
 significantly <y\_axis\_trend>  
 compared <y\_x\_comparison>  
 to  
 Europe <x\_axis\_label\_Scnd\_highest\_value>

*Computer science 60% and finally engineering at 70%*

Computer science <x\_axis\_label\_Scnd\_highest\_value>  
 60 <y\_axis\_Scnd\_highest\_val>  
 % <y\_axis\_inferred\_label>  
 and  
**finally** <order\_last>  
 engineering <x\_axis\_label\_highest\_value>  
 at  
 70 <y\_axis\_highest\_value\_val>  
 % <y\_axis\_inferred\_label>

## 27) y\_x\_comparison

Lexical expressions marking comparison are annotated with y\_x\_comparison. This applies to comparisons where the bar comparison is made by providing their heights, but also in cases where the heights are not mentioned explicitly. It can refer to a comparison of two or more bars.

The label annotates verbs ("compared"), nouns ("comparison"), conjunctions ("whilst") or any other part of speech that carries the same meaning.

*Insurance has the highest representation at 65% whereas law firms has the lowest.*

Insurance <x\_axis\_label\_highest\_value>  
has  
the  
highest <y\_axis\_highest\_value>  
representation <y\_axis>  
at  
65 <y\_axis\_highest\_value\_val>  
% <y\_axis\_inferred\_label>  
**whereas** <y\_x\_comparison>  
law firms <x\_axis\_label\_least\_value>  
has  
the  
lowest <y\_axis\_least\_value>

*32% spend time with family, while 38% prefer to read a book.*

32 <y\_axis\_Scnd\_highest\_val>  
% <y\_axis\_inferred\_label>  
spend time with family <x\_axis\_label\_Scnd\_highest\_value>  
,  
**while** <y\_x\_comparison>  
38 <y\_axis\_highest\_value\_val>  
% <y\_axis\_inferred\_label>  
prefer  
to  
read a book <x\_axis\_label\_highest\_value>

*This has reduced by over 50% compared to the year 2012.*

This  
has  
reduced <slope\_down>  
by

over <y\_axis\_approx>  
 50 <y\_axis\_inferred\_value\_mul\_v1=highest\_v2=least>  
 % <y\_axis\_inferred\_label>  
**compared** <y\_x\_comparison>  
 to  
 the  
 year <x\_axis>  
 2012 <x\_axis\_label\_highest\_value>

## 28) $t=*\_a=*$

The above listed labels sometimes do not suffice to annotate all chart- or topic-specific vocabulary. This vocabulary and its derivations can appear in the title or not, and modify the variable on y-, x- or both axes. We introduce the label  $t=[0|1]\_a=[x|y|b]$ , which should cover for such cases.

In the label

- "t" stands for "title" and is a binary marker of whether or not a word/phrase appears in the title (derivations included; synonyms not)
- "a" stands for axis; the value is either x, y or b for both axes.

Chart: Women representation in different sectors

- $t=1\_a=y$  : "tech has the lowest with 20% of women"
- $t=0\_a=y$ : "more women than men are found in law firms"

Chart: Median salary of software engineers per year wrt. their degree

- $t=1\_a=b$ : "63 k is median salary for software engineers who have PhD"
- $t=1\_a=y$ : "the median salary is at 55 k"

Chart: Gender pay gap

- $t=0\_a=x$ : "this chart show gender pay gap across 3 European countries"

Chart: Average time spent on social media daily in Maputo by age group

Sentence: In Maputo, persons of 45-54 spend around 150 minutes daily on social media.

- $t=1\_a=b$ : Maputo
- $t=1\_a=x$ : persons (fix label to  $t=0$ )
- $t=1\_a=y$ : daily; social media

## 29) other\_operation

Currently, the label set supports only pairwise comparisons, so all other cases are annotated with other\_operation

*Spain is exactly halfway between the two countries.*

Spain <x\_axis\_label\_Scnd\_highest\_value>  
is  
exactly  
**halfway** <other\_operation>  
between  
the  
two  
countries <x\_axis>

*STEM being the lowest at 15% and medicine being the highest at 31% with the other course of studies being in between them.*

STEM <x\_axis\_label\_least\_value>  
being  
the  
lowest <y\_axis\_least\_value>  
at  
15 <y\_axis\_least\_value\_val>  
% <y\_axis\_inferred\_label>  
and  
medicine <x\_axis\_label\_highest\_value>  
being  
the  
highest <y\_axis\_highest\_value>  
at  
31 <y\_axis\_highest\_value\_val>  
% <y\_axis\_inferred\_label>  
with  
the  
other  
course of studies <x\_axis\_labels\_rest>  
being  
in  
**between** <other\_operation>  
them

### 30) separator

In the preprocessing pipeline we wanted to keep the original structure of the text intact as much as possible. Newline characters often delimit prepositions, In the summaries, they are annotated with the separator label.

This  
chart  
shows  
the  
causes of obesity in Kiribati <topic>

.

\\n <separator>

The  
biggest <y\_axis\_highest\_value>  
cause <x\_axis>  
,  
shown  
to  
be  
responsible  
for  
40 <y\_axis\_inferred\_highest\_value\_approx>  
% <y\_axis\_inferred\_label>  
of  
cases <t=1\_a=y>  
,  
is  
genes <x\_axis\_label\_highest\_value>

### 31) interpretation

Some descriptions include a reading of the data that is not grounded in the chart or includes information that goes beyond what is accessible in the chart. In such cases, the <interpretation> label is used.

Some examples:

- However it doesn't show how this relates to inflation
- If this trend is to continue there were be a decrease in the following years
- This could be due to the economic situation in each respective continent, and also due to societal norms and values. In Asia , there is a lot of emphasis placed on the value of having a good education, which explains why more is spent on higher education in Asia than any other continent.

Note that inside the span of the interpretation label, we do not label entities as we normally would, e.g. "Asia" in the last example is left unlabeled.