# Recap day 4
## Robert Haase, Marcelo Zoccoler, Johannes Müller

October 2022

**PCA**: Linear transformation and identifies axes that explain most of the variance in the data

→ If the explained variances of $component_1$ and $component_2$ are [0.32, 0.2], we…

| Need more components | Cry | Will find no groups in the data |
|---|---|---|

→ If we add more data to a PCA-transformed set of data, we …
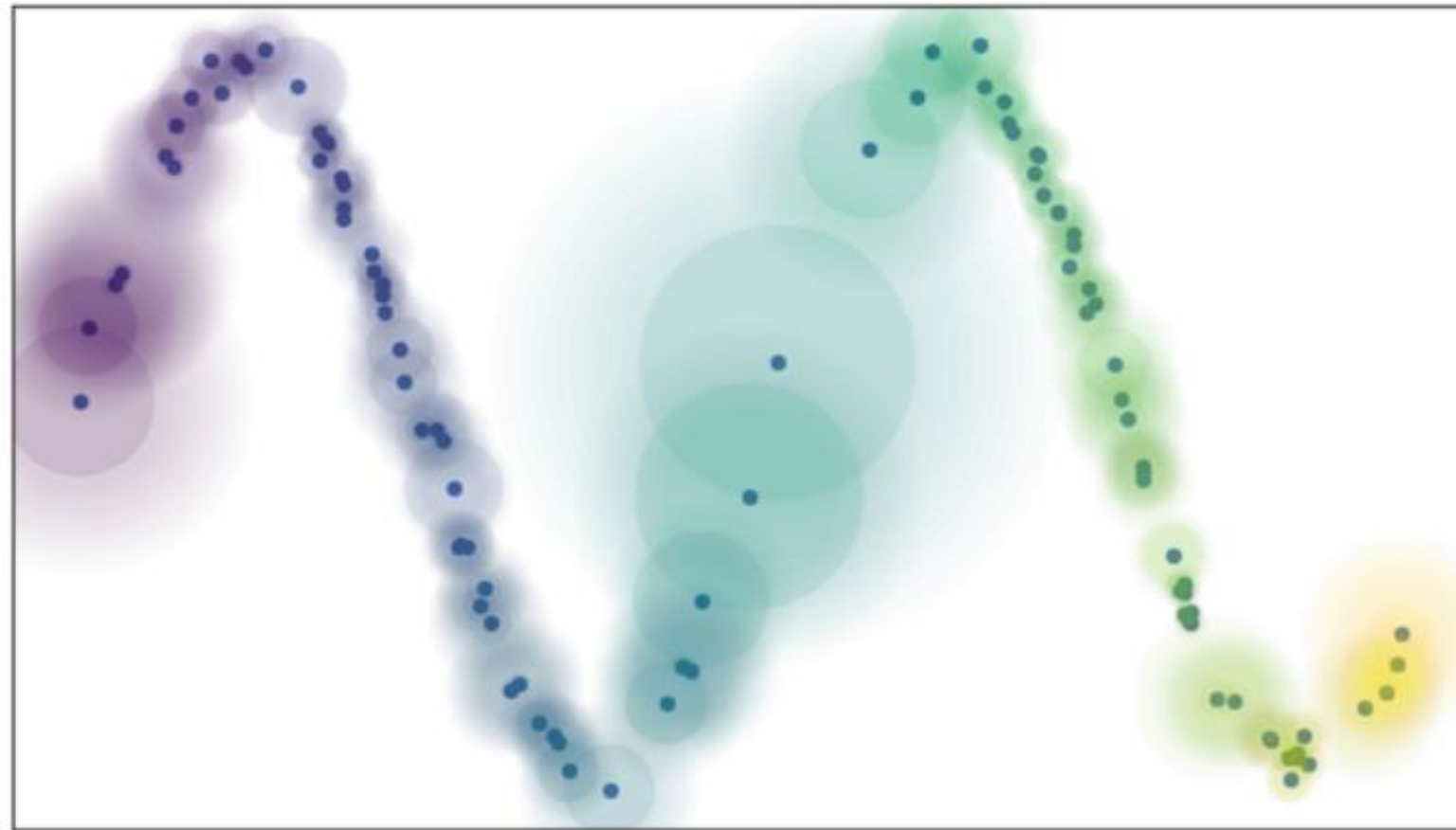
| Need to re-run the PCA | Can use the same PCA components | Cannot use PCA at all |
|---|---|---|

→ We can measure meaningful differences between data point in terms of $component_1$ and $component_2$

| True | False |
|---|---|

Explained variance $Component_1$: 0.98
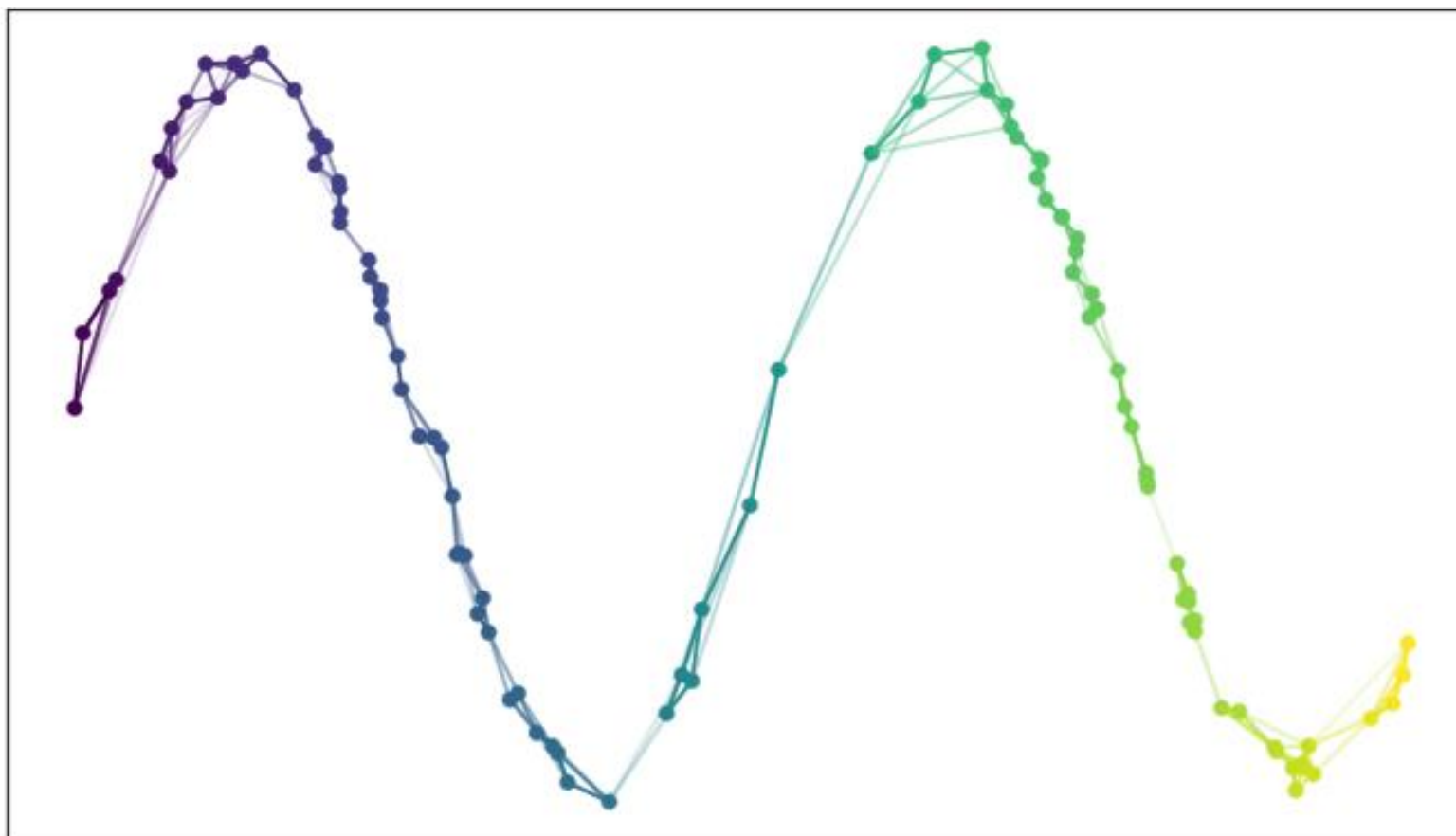Explained variance $Component_2$: 0.01

**UMAP:**
→ Link neighboring points by introducing point-wise distance metrics
→ "Relax" neighborhood graph to be able to display it in fewer dimensions while preserving its topology

We can measure meaningful differences between data points in terms of $UMAP_0$ and $UMAP_1$

| True | False |
|------|-------|

If we add new points to an existing UMAP projection, we can use the same projection

| True | False |
|------|-------|

Running a UMAP twice will give the same result

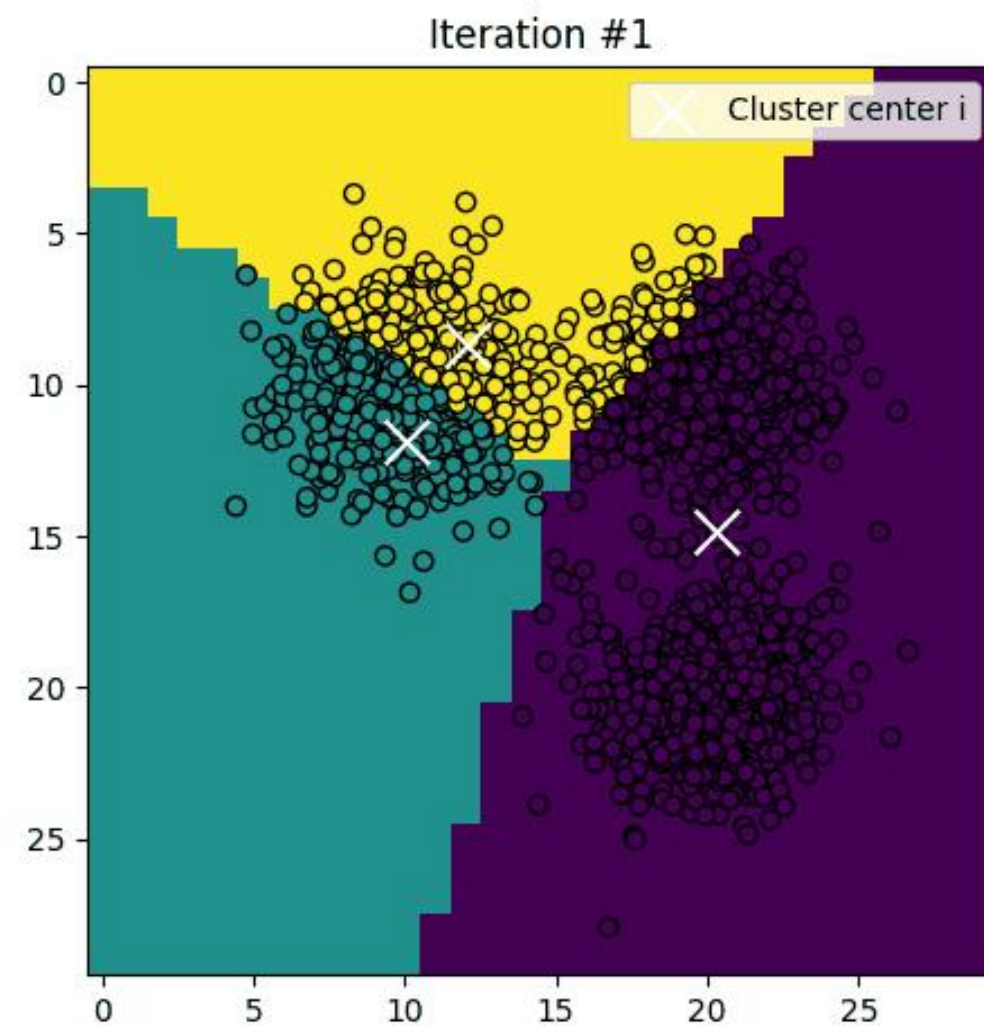| True | False |
|------|-------|



↑ 108 m · ↓ 0 m

221 m

113 m

Source: https://hdbscan.readthedocs.io/en/latest/how_hdbscan_works.html

# Clustering



**K-Means clustering:**

If we add more data to the pool, we can infer the cluster of a new point from the previously determined clusters

| True | False |
|------|-------|

**HDBSCAN:** If we add more data to the pool, we can infer the cluster of a new point from the previously determined clusters
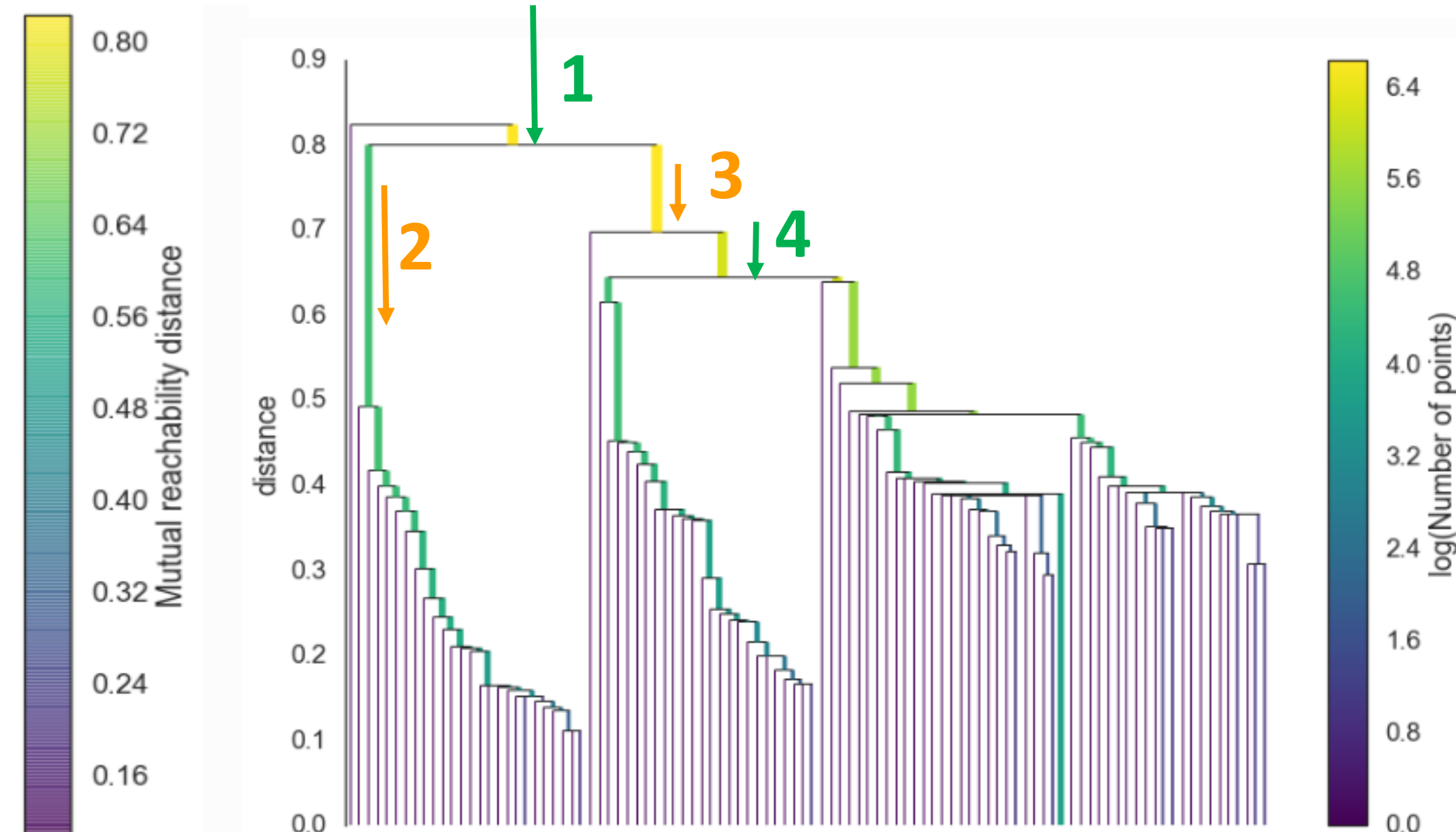


| True |
|------|
| False |

An important parameter is…

| The point density |
|-------------------|

| Data units |
|------------|

| Minimal cluster size |
|----------------------|

Source: https://hdbscan.readthedocs.io/en/latest/how_hdbscan_works.html

**October 2022**

## A Tutorial on Principal Component Analysis

[Jonathon Shlens](https://arxiv.org/abs/1404.1100)
[https://arxiv.org/abs/1404.1100](https://arxiv.org/abs/1404.1100)



$$X = \begin{bmatrix} x_1 \\ \vdots \\ x_m \end{bmatrix}$$

$C_X$ captures the covariance between all possible pairs of measurements. The covariance values reflect the noise and redundancy in our measurements.

- In the diagonal terms, by assumption, large values correspond to interesting structure.

- In the off-diagonal terms large magnitudes correspond to high redundancy.



*covariance matrix* $C_X$.

$$C_X \equiv \frac{1}{n}XX^T.$$

*Is there another basis, which is a linear combination of the original basis, that best re-expresses our data set?*

$$PX = Y$$

**SVD is a more general method for change of basis.**  $X = U\Sigma V^T$

$$
\begin{aligned}
C_Y &= \frac{1}{n}YY^T \\
&= \frac{1}{n}(PX)(PX)^T \\
&= \frac{1}{n}PXX^TP^T \\
&= P(\frac{1}{n}XX^T)P^T \\
C_Y &= PC_XP^T
\end{aligned}
$$

A symmetric matrix is diagonalized by a matrix of its orthonormal eigenvectors.

$$A = EDE^T$$

$$
\begin{aligned}
C_Y &= PC_XP^T \\
&= P(E^TDE)P^T
\end{aligned}
$$

**If we select P as the eigenvectors of $C_x$:** $P \equiv E^T$

$$
\begin{aligned}
&= P(P^TDP)P^T \\
&= (PP^T)D(PP^T) \\
&= (PP^{-1})D(PP^{-1}) \\
C_Y &= D
\end{aligned}
$$

# Quiz



## It's ok to reuse this picture on slides?

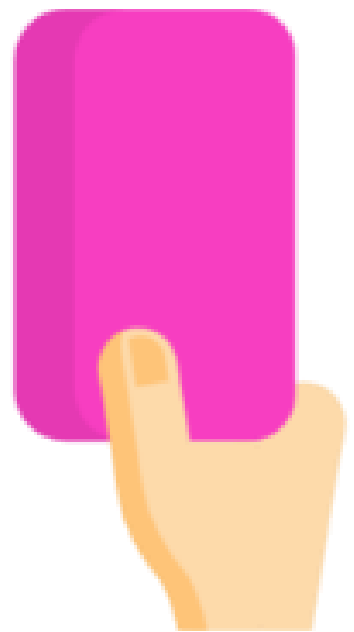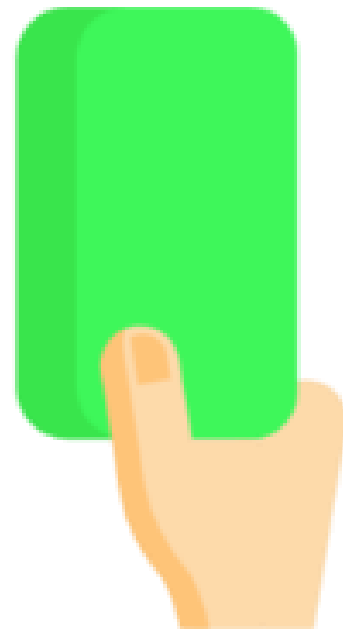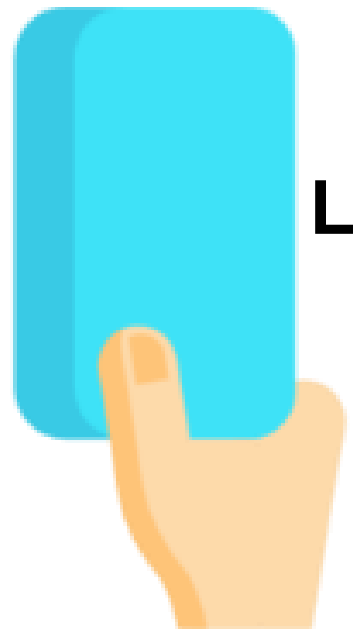Yes    No

# Quiz

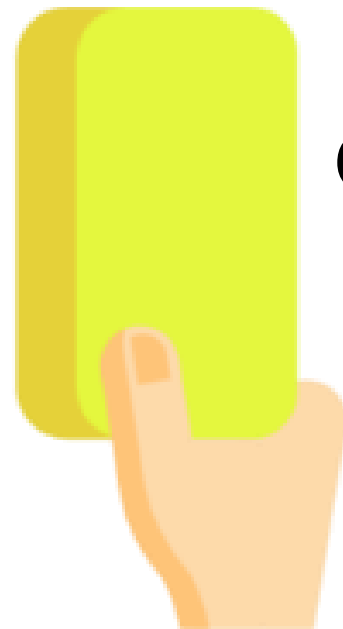It's ok to reuse code from this repository if...
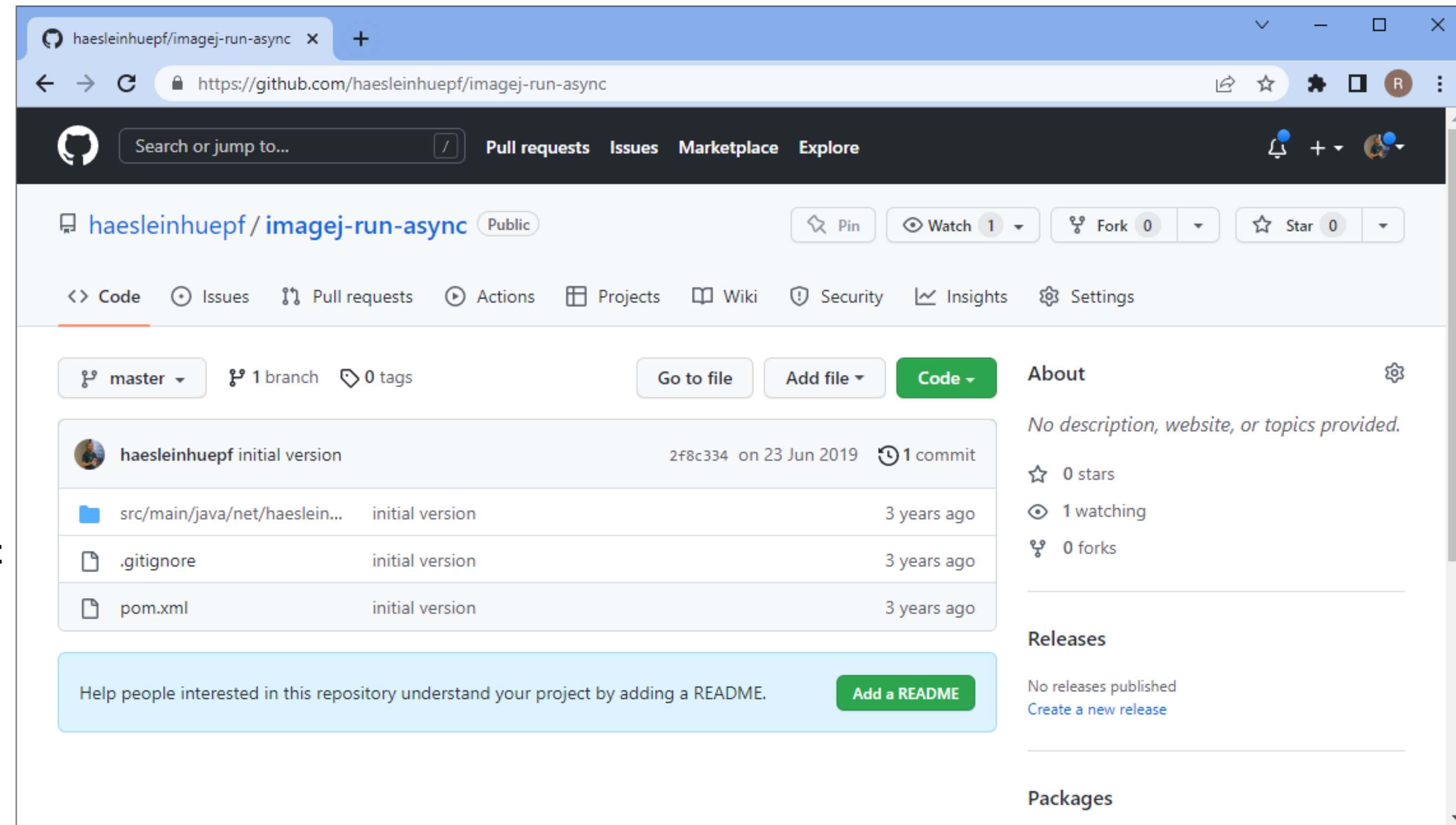
**Mention author**

**Ask the authors**
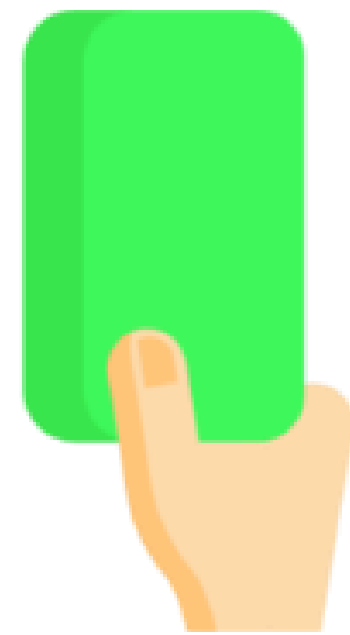
**Link to the repository**

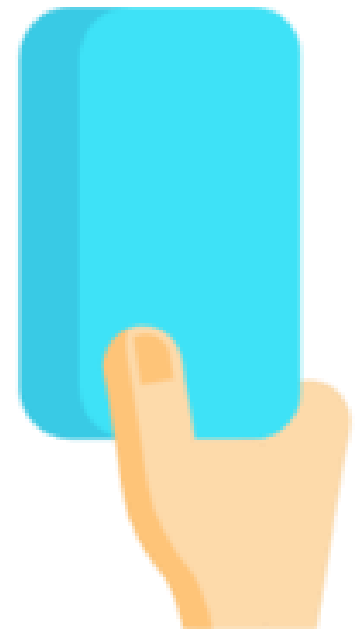**Copy the copyright statement**

# Quiz

## It's ok to reuse Figures from XKCD.com if…

**Mention author**
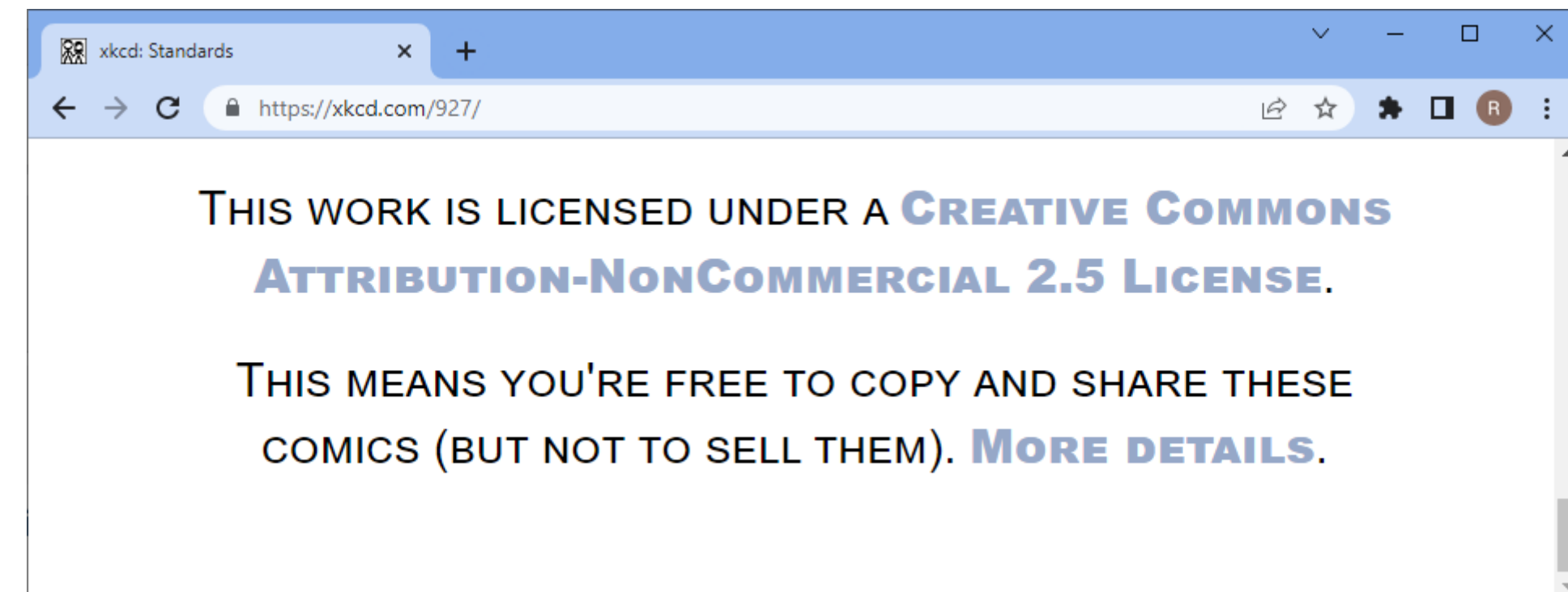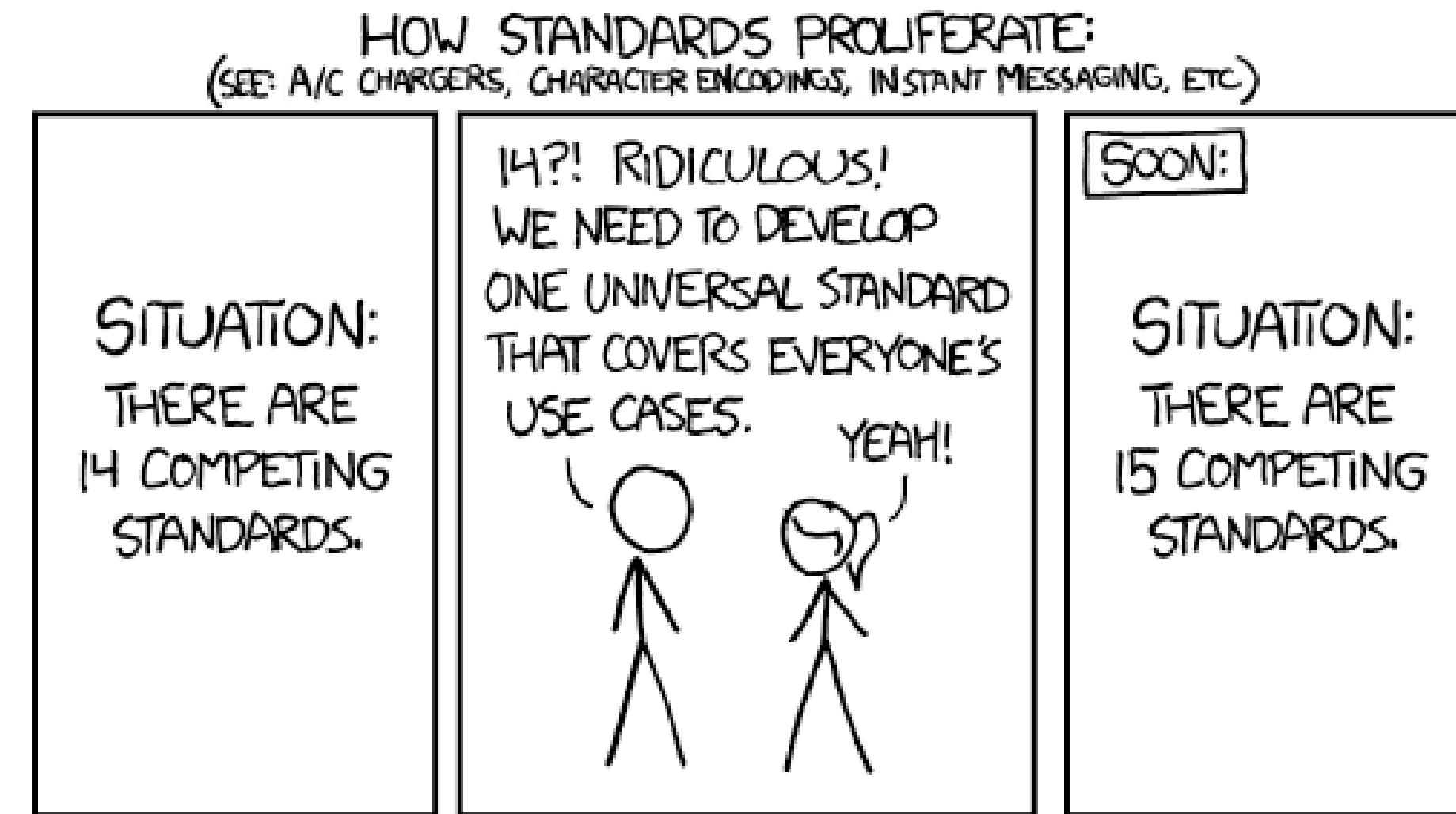
**Ask the authors**

**Link to the license**

**Copy the copyright statement**



https://xkcd.com/927/
The Figure is licensed by XKCD.com under a **Creative Commons Attribution-NonCommercial 2.5 License**.