# Code Duplication and Reuse in Jupyter Notebooks

**Andreas P. Koenzen, Neil A. Ernst, Margaret-Anne D. Storey**

akoenzen@uvic.ca, nernst@uvic.ca, mstorey@uvic.ca

University of Victoria

Victoria, Canada

**Presenter: Andreas P. Koenzen**

# Studies Conducted

**Study #1:** Quantified self-duplicated code snippets in Jupyter notebooks at the repository level **(Artifacts)**

**Study #2:** Observed participants solving tasks using Jupyter notebooks **(Behaviour)**

✓ **Motivation**　　　✓ Study #1　　　✓ Study #2　　　✓ Limitations　　　✓ Future Work　　　✓ Conclusion

**Presenter: Andreas P. Koenzen (University of Victoria, Canada)**　　　Code Duplication and Reuse in Jupyter Notebooks - VL/HCC 2020 (August 2020)

# Why..?

**Does it happen..? + What method..? + From where..? = Understanding code reuse can lead to tools that expedite data exploration using Jupyter notebooks**

✓ Motivation    ✓ Study #1    ✓ Study #2    ✓ Limitations    ✓ Future Work    ✓ Conclusion

Presenter: Andreas P. Koenzen (University of Victoria, Canada)    Code Duplication and Reuse in Jupyter Notebooks - VL/HCC 2020 (August 2020)

# Code Duplication (Artifact)

**RQ1:** How much cell code duplication occurs in Jupyter Notebooks?

**Artifacts**



```
plt.figure(figsize=(7,7),dpi=400)    plt.figure(figsize=(7,7),dpi=400)
ax = plt.subplot(2,1,1)              ax = plt.subplot(2,1,1)
plot(PPT, Nash_Flow, 'bo',           plot(H1, r2_Sed, 'bo')
markersize=3)
title('PPT', fontsize=14.,           title('H1', fontsize=14.,
y=1.02, fontweight='bold')           y=1.02, fontweight='bold')
ax = plt.subplot(2,1,2)              ax = plt.subplot(2,1,2)
plt.hist(PPT)                        plt.hist(H1)
np.corrcoef(PPT,Nash_Flow)           np.corrcoef(H1,r2_Sed)
```

Fig. 4: Example of a Type-2 duplicate detected by our algorithm with Levenshtein distance of 42 and Duplicate Ratio of 0.27. Coded as *Visualization*.

✓ Motivation     ✓ Study #1     ✓ Study #2     ✓ Limitations     ✓ Future Work     ✓ Conclusion

Presenter: Andreas P. Koenzen (University of Victoria, Canada)     Code Duplication and Reuse in Jupyter Notebooks - VL/HCC 2020 (August 2020)

# Code Reuse (Behaviour)

**RQ2:** How does cell code reuse happen in Jupyter notebooks?

**How these artifacts are introduced into notebooks**

**Method**

# Code Reuse (Behaviour)

**RQ3:** What are the preferred sources for code reuse in Jupyter notebooks?

**From where these artifacts are introduced into notebooks**

**Source**



Image source: https://www.vecteezy.com/free-vector/www-icon

✓ Motivation          ✓ Study #1          ✓ Study #2          ✓ Limitations          ✓ Future Work          ✓ Conclusion

Presenter: Andreas P. Koenzen (University of Victoria, Canada)          Code Duplication and Reuse in Jupyter Notebooks - VL/HCC 2020 (August 2020)
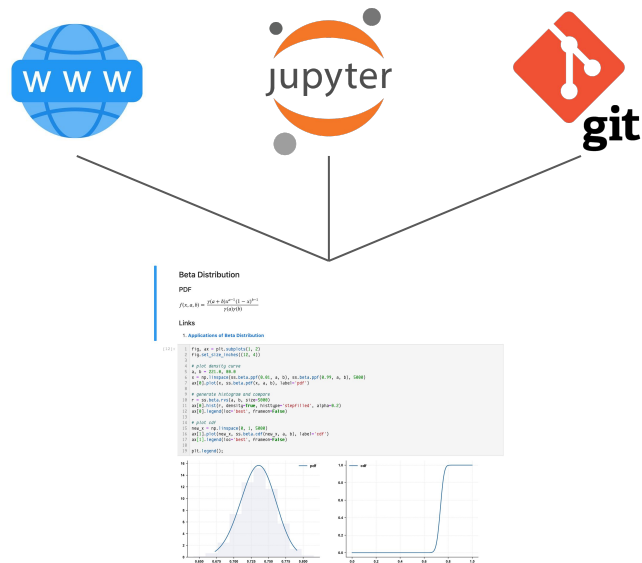
# Study #1

- Analyzed GitHub repositories that contained at least 1 Jupyter notebook and computing clones within these repositories

- Thematically coded clones detected before, to assign a computational purpose

# Study #1 / Methodology

- Randomly sampled 1,000 repositories from GitHub

- Computed duplicates (Type-1, Type-2 and Type-3) with custom function

- Randomly sampled 500 detected clones from previous step and thematically coded them based on the programming task these clones performed e.g. data visualization
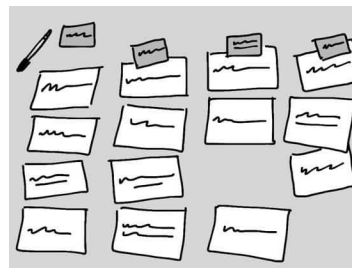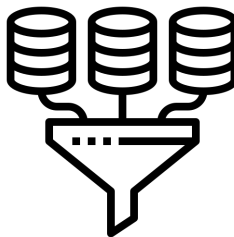
Image source: https://www.flaticon.com/authors/becris

# Study #1 / Results

- 1 in 13 code cells are clones
- Visualization routines are duplicated the most

✓ Motivation     ✓ Study #1     ✓ Study #2     ✓ Limitations     ✓ Future Work     ✓ Conclusion

Presenter: Andreas P. Koenzen (University of Victoria, Canada)     Code Duplication and Reuse in Jupyter Notebooks - VL/HCC 2020 (August 2020)

# Study #2

- Observed users perform previously designed specific tasks using Jupyter notebooks

Conducted at the CHISEL lab in the University of Victoria with 8 participants (All students)

# Study #2 / Methodology

- Observational study in a lab setting
- Participants had to solve 3 tasks using Jupyter notebooks according to self-reported expertise
- Minimal limitations (Extra time was given if necessary)
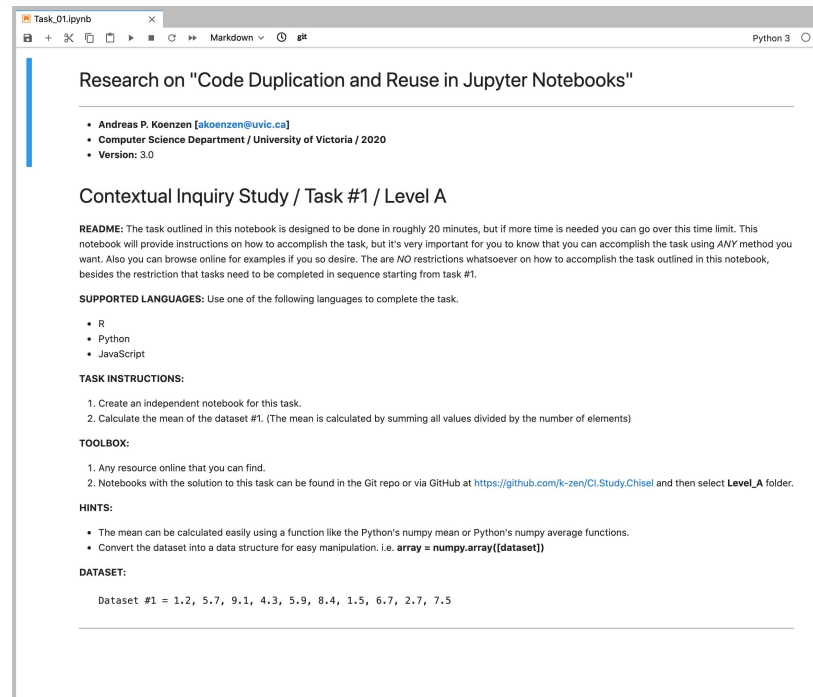- Solutions to tasks were available in the study's git repository



Task_01.ipynb

Markdown    git    Python 3

## Research on "Code Duplication and Reuse in Jupyter Notebooks"

- **Andreas P. Koenzen [akoenzen@uvic.ca]**
- **Computer Science Department / University of Victoria / 2020**
- **Version:** 3.0

### Contextual Inquiry Study / Task #1 / Level A

**README:** The task outlined in this notebook is designed to be done in roughly 20 minutes, but if more time is needed you can go over this time limit. This notebook will provide instructions on how to accomplish the task, but it's very important for you to know that you can accomplish the task using *ANY* method you want. Also you can browse online for examples if you so desire. The are *NO* restrictions whatsoever on how to accomplish the task outlined in this notebook, besides the restriction that tasks need to be completed in sequence starting from task #1.

**SUPPORTED LANGUAGES:** Use one of the following languages to complete the task.

- R
- Python
- JavaScript

**TASK INSTRUCTIONS:**

1. Create an independent notebook for this task.
2. Calculate the mean of the dataset #1. (The mean is calculated by summing all values divided by the number of elements)

**TOOLBOX:**

1. Any resource online that you can find.
2. Notebooks with the solution to this task can be found in the Git repo or via GitHub at https://github.com/k-zen/CI.Study.Chisel and then select **Level_A** folder.

**HINTS:**

- The mean can be calculated easily using a function like the Python's numpy mean or Python's numpy average functions.
- Convert the dataset into a data structure for easy manipulation. i.e. **array = numpy.array([dataset])**

**DATASET:**

    Dataset #1 = 1.2, 5.7, 9.1, 4.3, 5.9, 8.4, 1.5, 6.7, 2.7, 7.5

✓ Motivation    ✓ Study #1    ✓ Study #2    ✓ Limitations    ✓ Future Work    ✓ Conclusion

Presenter: Andreas P. Koenzen (University of Victoria, Canada)    Code Duplication and Reuse in Jupyter Notebooks - VL/HCC 2020 (August 2020)

# Study #2 / Methodology

- Video coding
- Audio transcriptions
- Notes
- Questions based on observations. i.e. why did you used that particular method
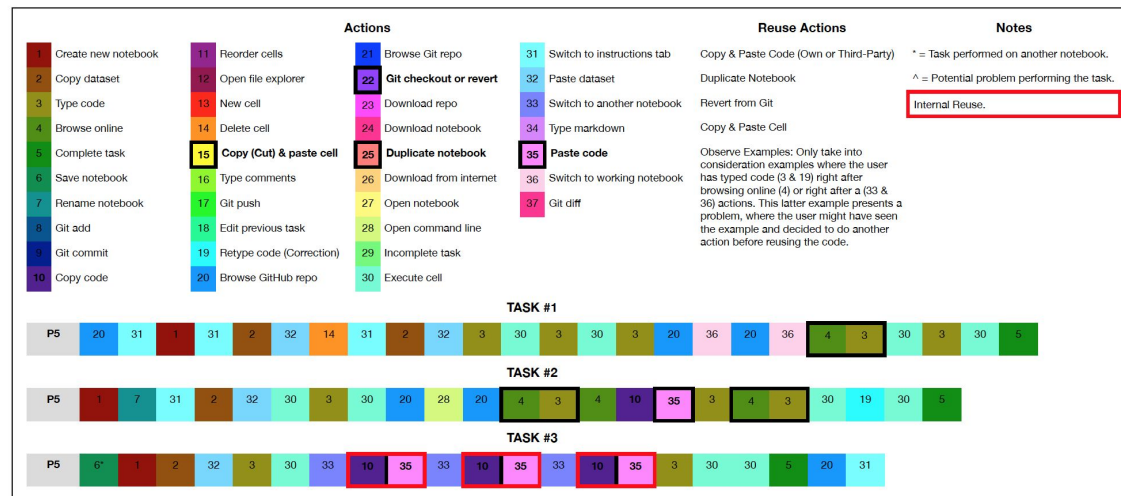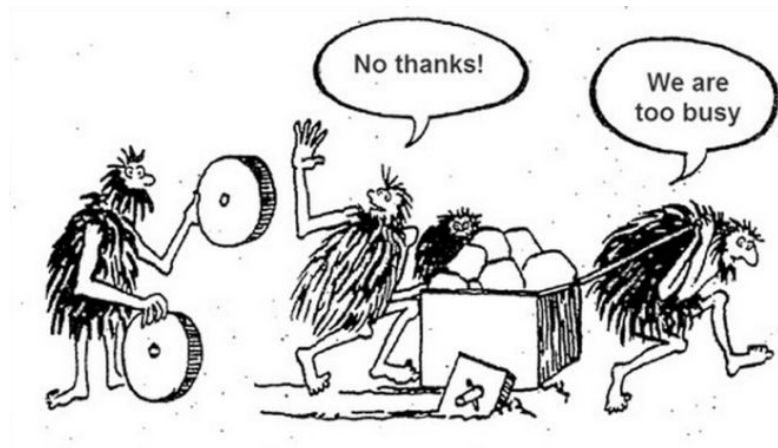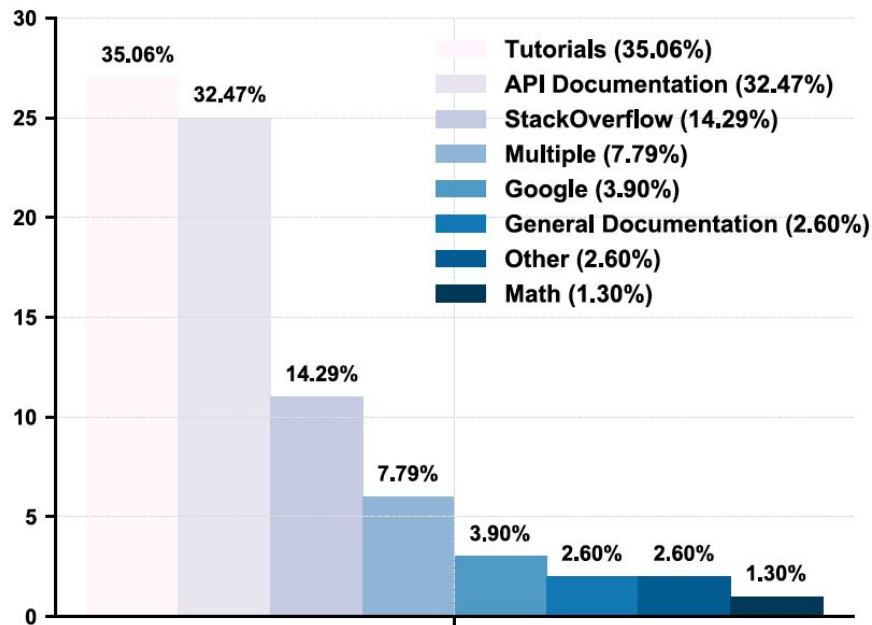- Short interview and questionnaire



Fig. 2: Example coding of steps one of our participants (P5) made during the observational study, based on video and audio recordings.

✓ Motivation          ✓ Study #1          ✓ Study #2          ✓ Limitations          ✓ Future Work          ✓ Conclusion

Presenter: Andreas P. Koenzen (University of Victoria, Canada)          Code Duplication and Reuse in Jupyter Notebooks - VL/HCC 2020 (August 2020)

# Study #2 / Results

- Participants reused extensively
- External libraries like "numpy" were used extensively
- Most common source of reuse: online sources (18% of total time)
- Least common source of reuse: VCS (0 participants)
- Reutilization was done through copying and pasting, copying by typing of code
- Least common method of reuse: duplicating a notebook

✓ Motivation     ✓ Study #1     ✓ Study #2     ✓ Limitations     ✓ Future Work     ✓ Conclusion

Presenter: Andreas P. Koenzen (University of Victoria, Canada)     Code Duplication and Reuse in Jupyter Notebooks - VL/HCC 2020 (August 2020)

# Study #2 / Results



| Code | Overall Count | Task #1 | Task #2 | Task #3 |
|---|---|---|---|---|
| C&P | 20 times | 0 times | 8 times | 12 times |
| CELL | 1 times | 0 times | 1 times | 0 times |
| TYPE | 0 times | 0 times | 0 times | 0 times |
| DUPE | 0 times | 0 times | 0 times | 0 times |
| GIT | 0 times | 0 times | 0 times | 0 times |
| TYPE_ON | 36 times | 16 times | 14 times | 6 times |
| NONE | 1 times | 0 times | 0 times | 1 times |

**Table 4.2:** Count of reuse codes for all participants and across all tasks. Highlighted in red are the highest counts.

✓ Motivation          ✓ Study #1          ✓ Study #2          ✓ Limitations          ✓ Future Work          ✓ Conclusion

Presenter: Andreas P. Koenzen (University of Victoria, Canada)          Code Duplication and Reuse in Jupyter Notebooks - VL/HCC 2020 (August 2020)

# Limitations

**Construct Validity**

Function parameters were found empirically (Cut-off value, λ) => Grid search using an oracled data set

Participants were constrained to use lab equipment => Instead of using their own

**Internal Validity**

Self-assessed level of proficiency => Might introduce Observer-Expectancy bias

Tasks might have been too simple => But we asked proficiency before 👆

**External Validity**

Notebooks were sampled from GitHub, which may differ from corporate settings

Students differ from industry practitioners

**Our findings should be seen as restricted to the sample we used**

✓ Motivation    ✓ Study #1    ✓ Study #2    ✓ Limitations    ✓ Future Work    ✓ Conclusion

Presenter: Andreas P. Koenzen (University of Victoria, Canada)    Code Duplication and Reuse in Jupyter Notebooks - VL/HCC 2020 (August 2020)

# Future Work

- Observational studies in real settings where participants used their own tools to solve real problems (Advanced students, industry practitioners)

- Harder and more complex problems might shed light into reutilization of complex routines via methods not observed during this study

✓ Motivation          ✓ Study #1          ✓ Study #2          ✓ Limitations          ✓ Future Work          ✓ Conclusion

Presenter: Andreas P. Koenzen (University of Victoria, Canada)          Code Duplication and Reuse in Jupyter Notebooks - VL/HCC 2020 (August 2020)

# Summary

- 1 in 13 code cells are clones
- Visualization routines are duplicated the most
- Users reuse extensively when using Jupyter notebooks
- The most common source of reuse is the web
- git only for storage
- Less reuse from own code (Reinvent the wheel attitude)

# So what...

- Code reuse support via tools like **Google Colab's Code Snippets** can be beneficial and might speed up the development process
- Extensive codebase should be put into modules (JupyterLab's **"autoreload"** magic)
- **Light version control** with a simple interface is more suitable for Jupyter notebooks
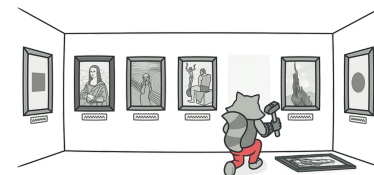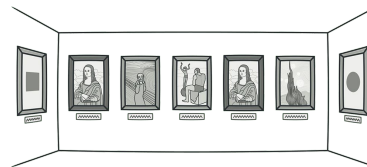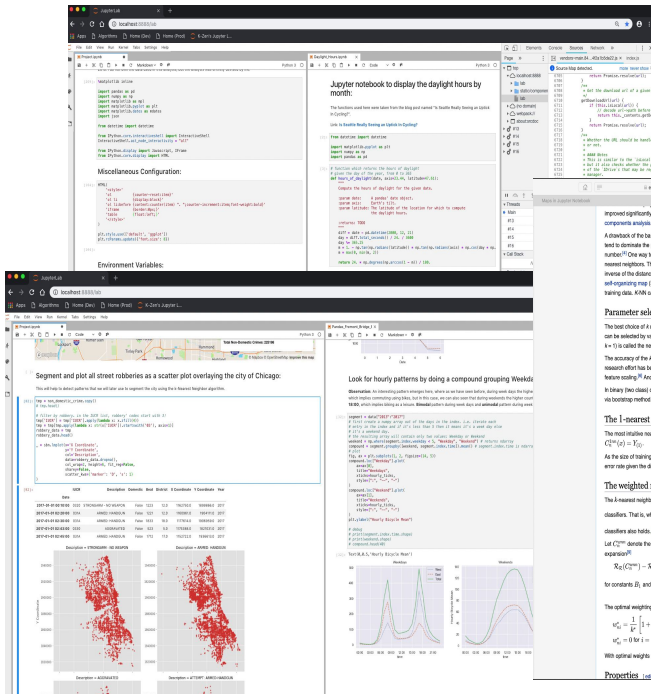
Image source: https://refactoring.guru/smells/duplicate-code

✓ Motivation          ✓ Study #1          ✓ Study #2          ✓ Limitations          ✓ Future Work          ✓ Conclusion

Presenter: Andreas P. Koenzen (University of Victoria, Canada)          Code Duplication and Reuse in Jupyter Notebooks - VL/HCC 2020 (August 2020)

# My own personal experience…

✓ Motivation    ✓ Study #1    ✓ Study #2    ✓ Limitations    ✓ Future Work    ✓ Conclusion

Presenter: Andreas P. Koenzen (University of Victoria, Canada)    Code Duplication and Reuse in Jupyter Notebooks - VL/HCC 2020 (August 2020)

# Q&A



**Reminder:**
- 1 in 13 code cells are clones
- Visualization routines are duplicated the most
- Users reuse extensively when using Jupyter notebooks
- The most common source of reuse is the web
- git only for storage
- Less reuse from own code (Reinvent the wheel attitude)

**Andreas P. Koenzen, Neil A. Ernst, Margaret-Anne D. Storey**
akoenzen@uvic.ca, nernst@uvic.ca, mstorey@uvic.ca
University of Victoria
Victoria, Canada
Pre-print: https://arxiv.org/abs/2005.13709