# Analysis of the Linux Kernel Evolution Using Code Clone Coverage

Simone Livieri

Yoshiki Higo
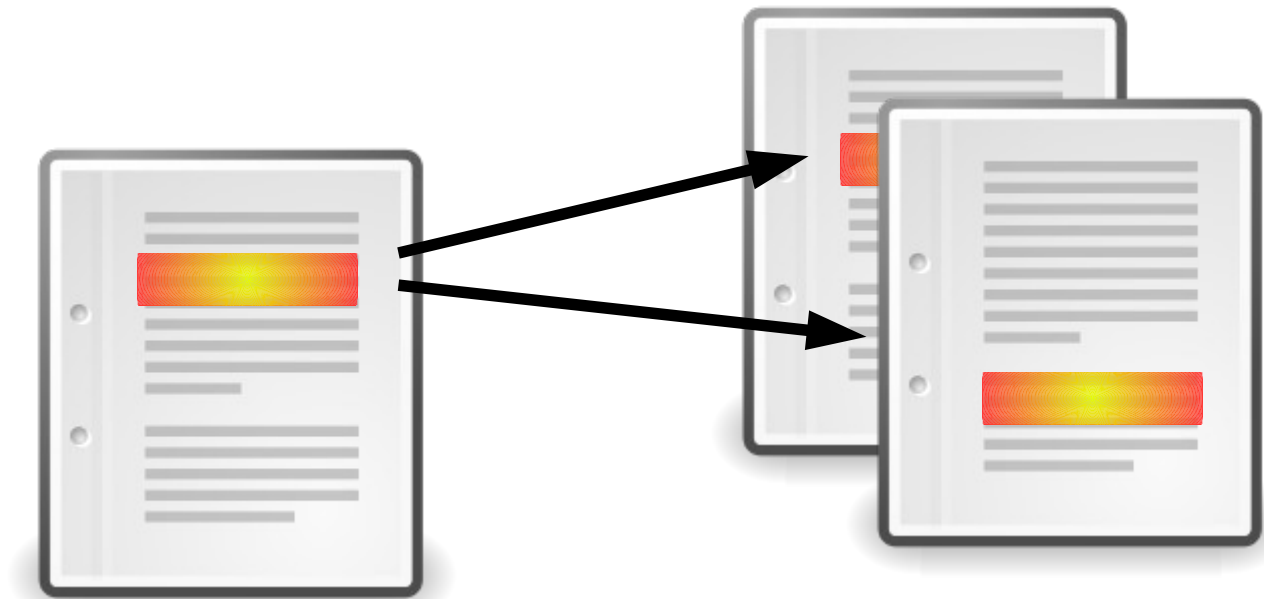
Makoto Matsushita

Katsuro Inoue

Software
Engineering
Laboratory

**Department of Computer Science, Graduate School of Information Science and Technology, Osaka University**

# Code clone?

# Code Clone?

- A code clone is a set of identical or similar fragments of code

# Code Clone Detection

- Various detection methods

  - Token based

  - Abstract Syntax Tree based

  - Program Dependence Graph based

  - …

Software
Engineering
Laboratory

# CCFinder

- Token-based code-clone detection tool

  - Insensitive to renamed variable and code layout

  - Multi language support (C, C++, COBOL, Java, …)

- Very good scalability and speed…

# CCFinder

- Token-based code-clone detection tool
  - Insensitive to renamed variable and code layout
  - Multi language support (C, C++, COBOL, Java, …)

- Very good scalability and speed…

- …but scalability is limited by the hardware used

Software
Engineering
Laboratory

# D-CCFinder

- A system for distributed code clone analysis

- Uses CCFinder as code clone detector

# What?

# What

Start a large scale study of a software system's evolution using code clone analysis

# Why?

Software
Engineering
Laboratory

# Why

- The evolution of a software system can be reconstructed with code clone analysis

- No large scale study have been performed yet

- D-CCFinder permits large scale code clone analysis

Software
Engineering
Laboratory

# Why

- The evolution of a software system can be reconstructed with code clone analysis

- No large scale study have been performed yet

- D-CCFinder permits large scale code clone analysis

- There were two weeks left before the deadline and we had nothing to do

# The guinea pig

Software
Engineering
Laboratory

# The Linux Kernel

- 15 years long development effort involving hundreds of developers

- Two development branches: **stable** and **unstable** (prior to version 2.6)

- The source code size grew from 3.8Mbytes (version 1.0) to 157Mbytes (2.6.18.3)

Software
Engineering
Laboratory

# The Linux Kernel

| Version | LOC | Size (Kbytes) | # of versions |
|---|---|---|---|
| 1.0 | 141K | 3,926 | 1 |
| 1.2.0~ | 234K | 6,534 | 14 |
| 1.2.13 | 238K | 6,596 | |
| 2.0.0~ | 563K | 16,076 | 41 |
| 2.0.40 | 768K | 21,952 | |
| 2.2.0~ | 1,310K | 37,056 | 27 |
| 2.2.26 | 1,970K | 58,812 | |
| 2.4.0~ | 2,366K | 69,200 | 34 |
| 2.4.33.4 | 3,865K | 112,148 | |
| 2.6.0~ | 4,120K | 120,030 | 19 |
| 2.6.18.3 | 5,476K | 157,290 | |

| | |
|---|---|
| **Total number of versions** | 136 |
| **Number of .c files** | 376,596 |
| **Total lines of code** | 266,943,565 |
| **Total size** | 7.4 Gbytes |

- 136 kernel versions from the stable branches

- Considered only .c files

- Size measured with `du`

- LOC counted with `wc`

Software Engineering Laboratory

# How?

# How

- For each pair (A, B) of kernel versions we computed and plotted the code clone coverage

$$Coverage(A,B) = \frac{Loc(CC(A,B))}{Loc(A) + Loc(B)}$$

CC(A,B)= code clone fragments between A and B
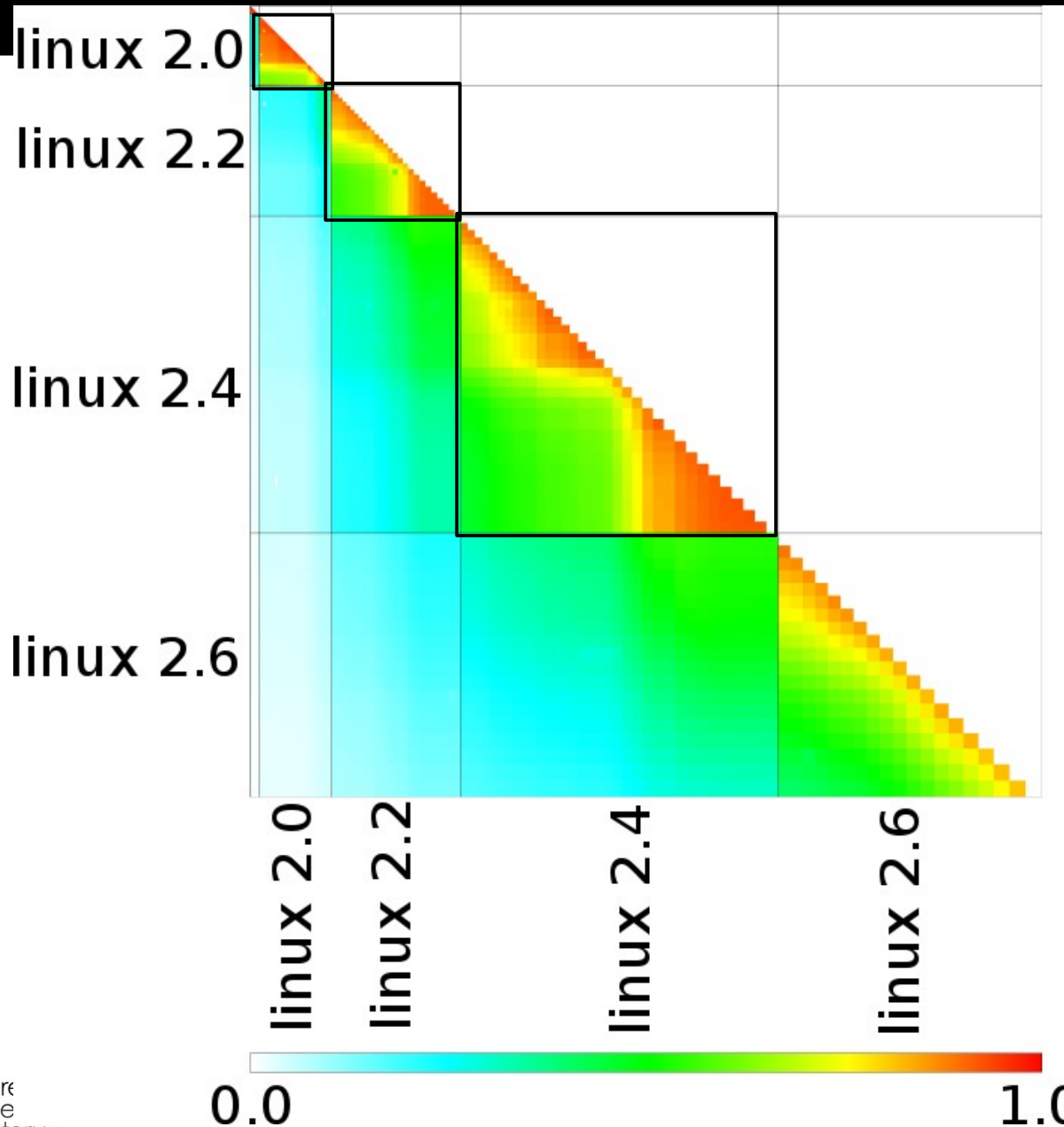
Software
Engineering
Laboratory

# Results

# Results

# Results



- Max coverage: 67%
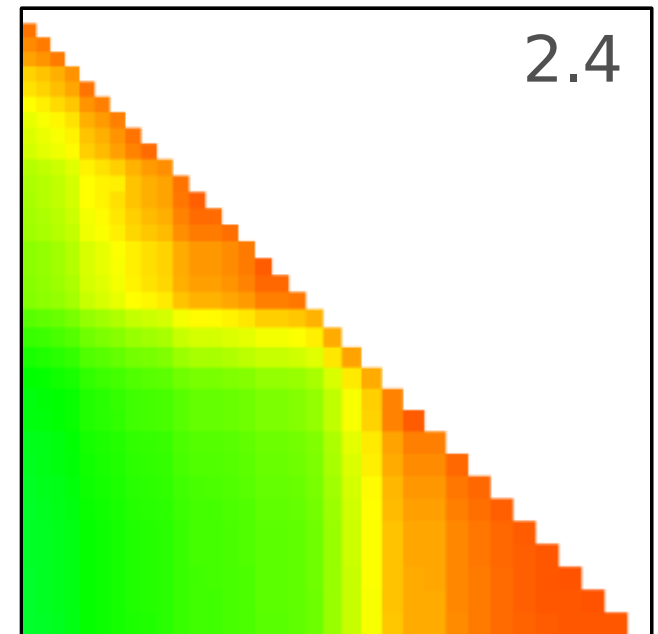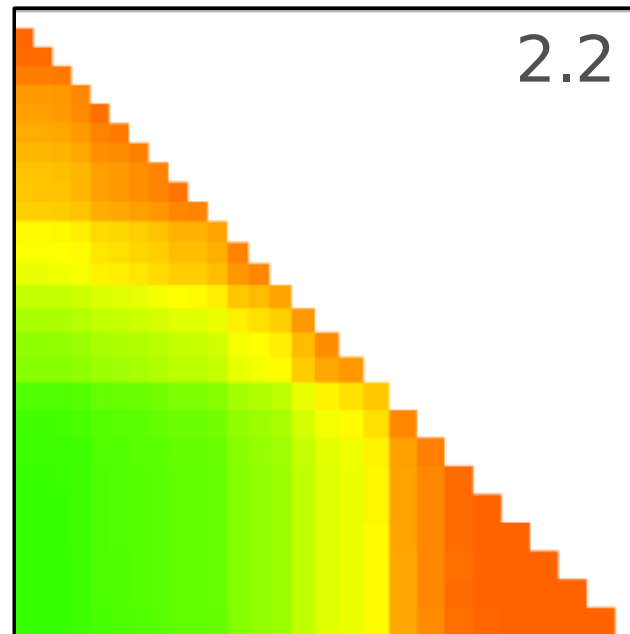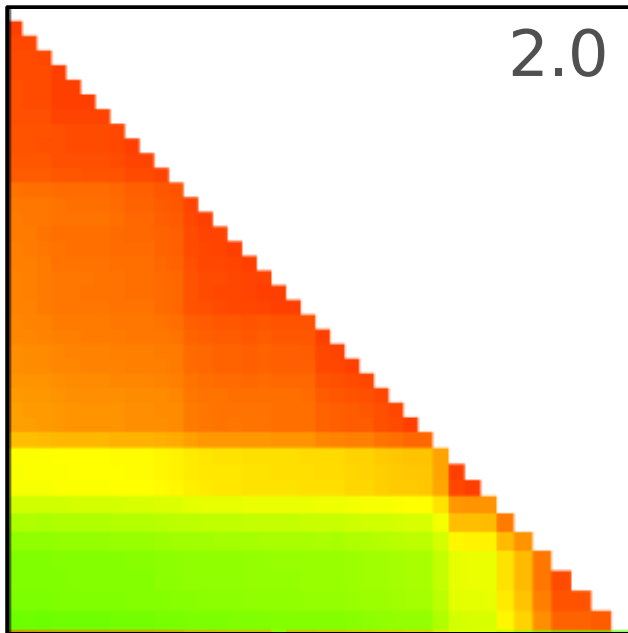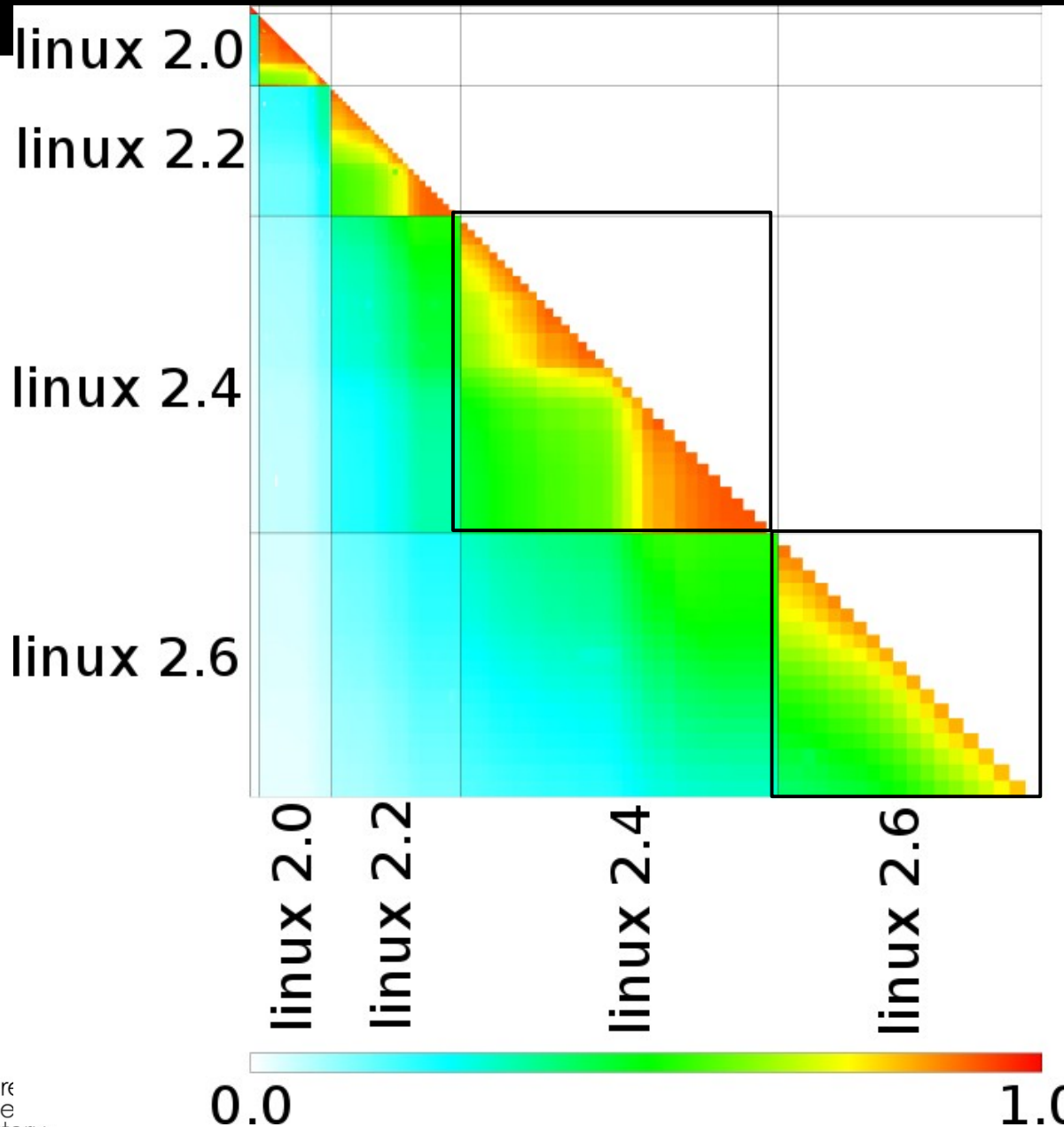- Highest coverage on the diagonal (as expected)

# Results

# Results

- Same pattern

- Code "back-ported" from the development branches



2.0

2.2

2.4

Software
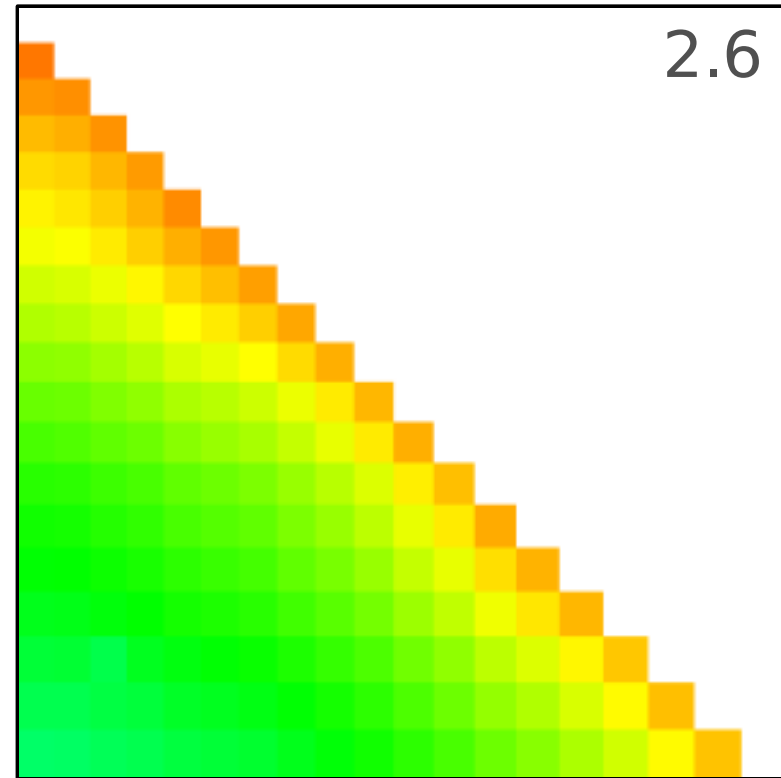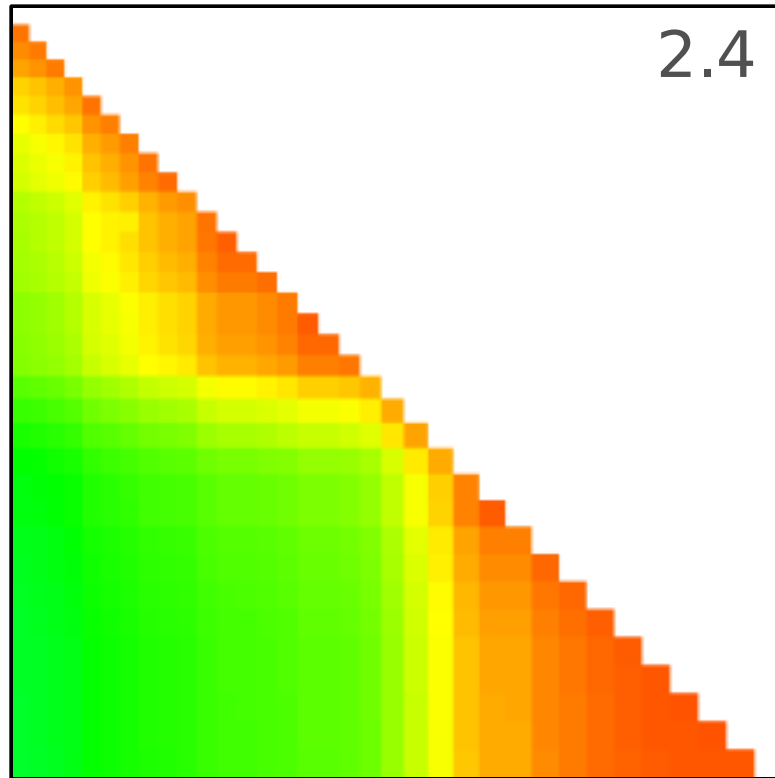Engineering
Laboratory

# Results

# Results

- Different patterns due to different development processes

# Conclusion & Future

- ## Conclusion

  - Presented a tentative study of the evolution of the Linux kernel computing and visualizing the code-clone coverage metric

- ## Future

  - Elaborate and complete the analysis presented and produce a more detailed report of the code changes

以上