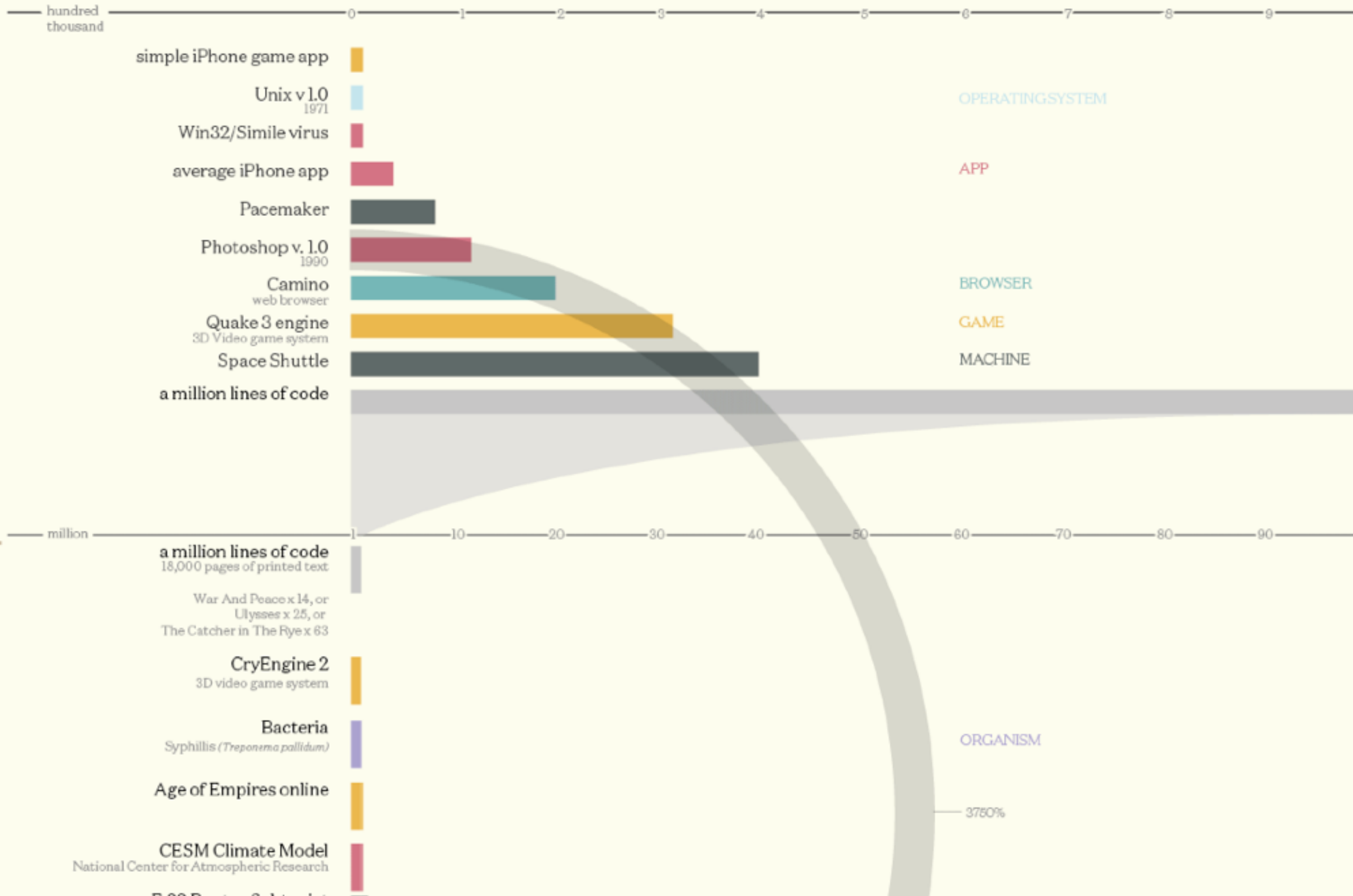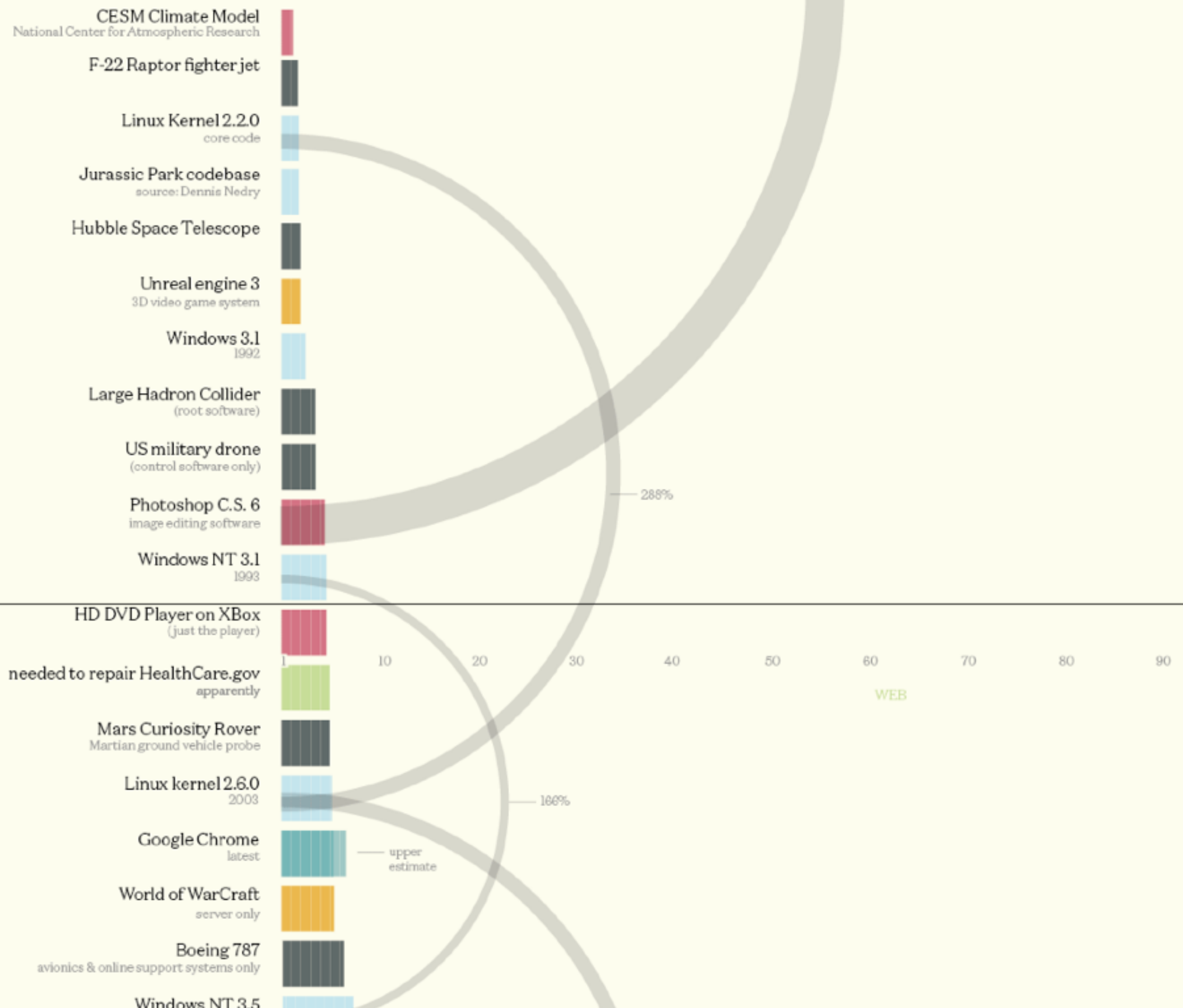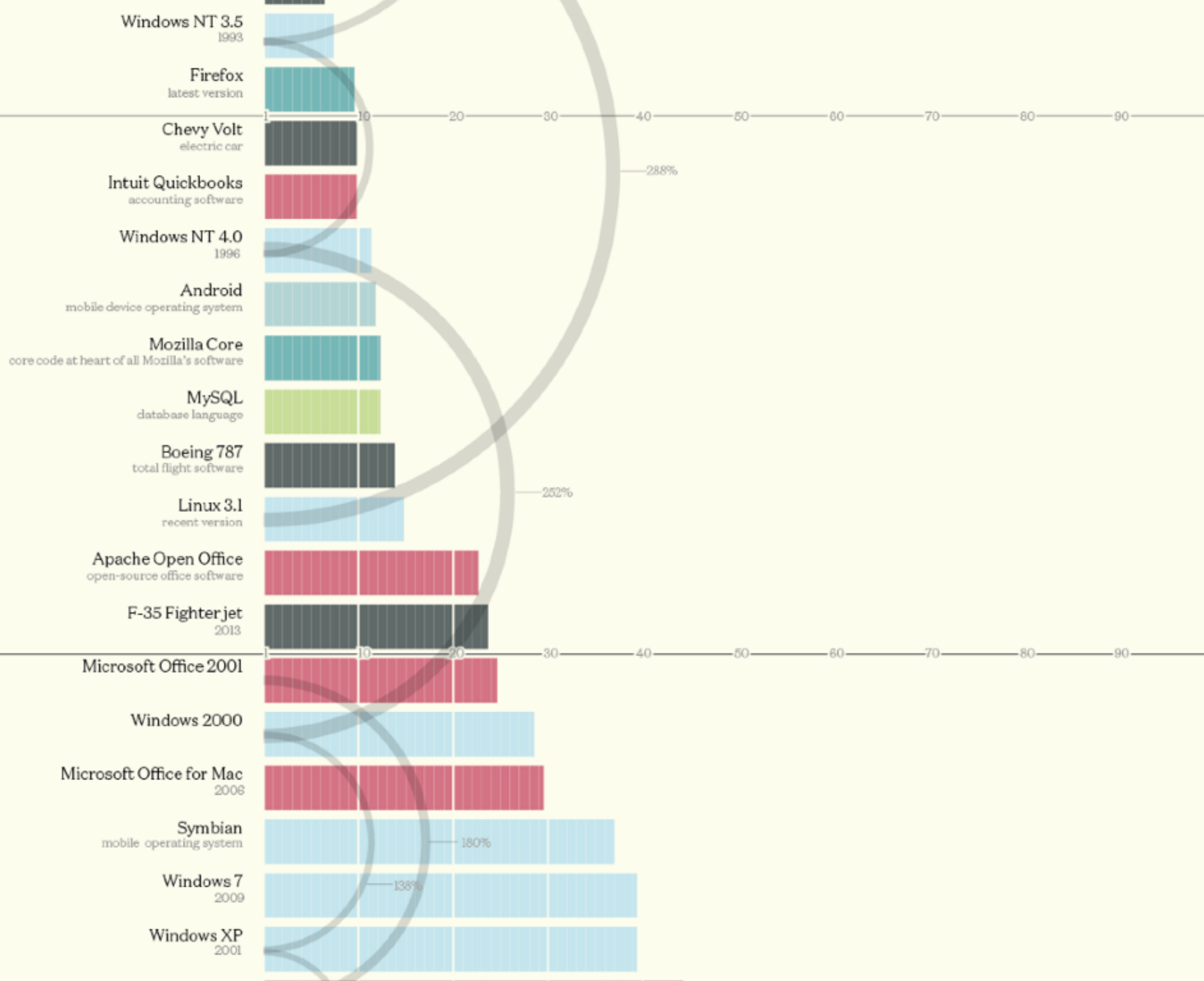# Introduction to Software Evolution

Dr. Vadim Zaytsev aka @grammarware
UvA, MSc SE, 25 October 2015

# Codebases

Millions of lines of code

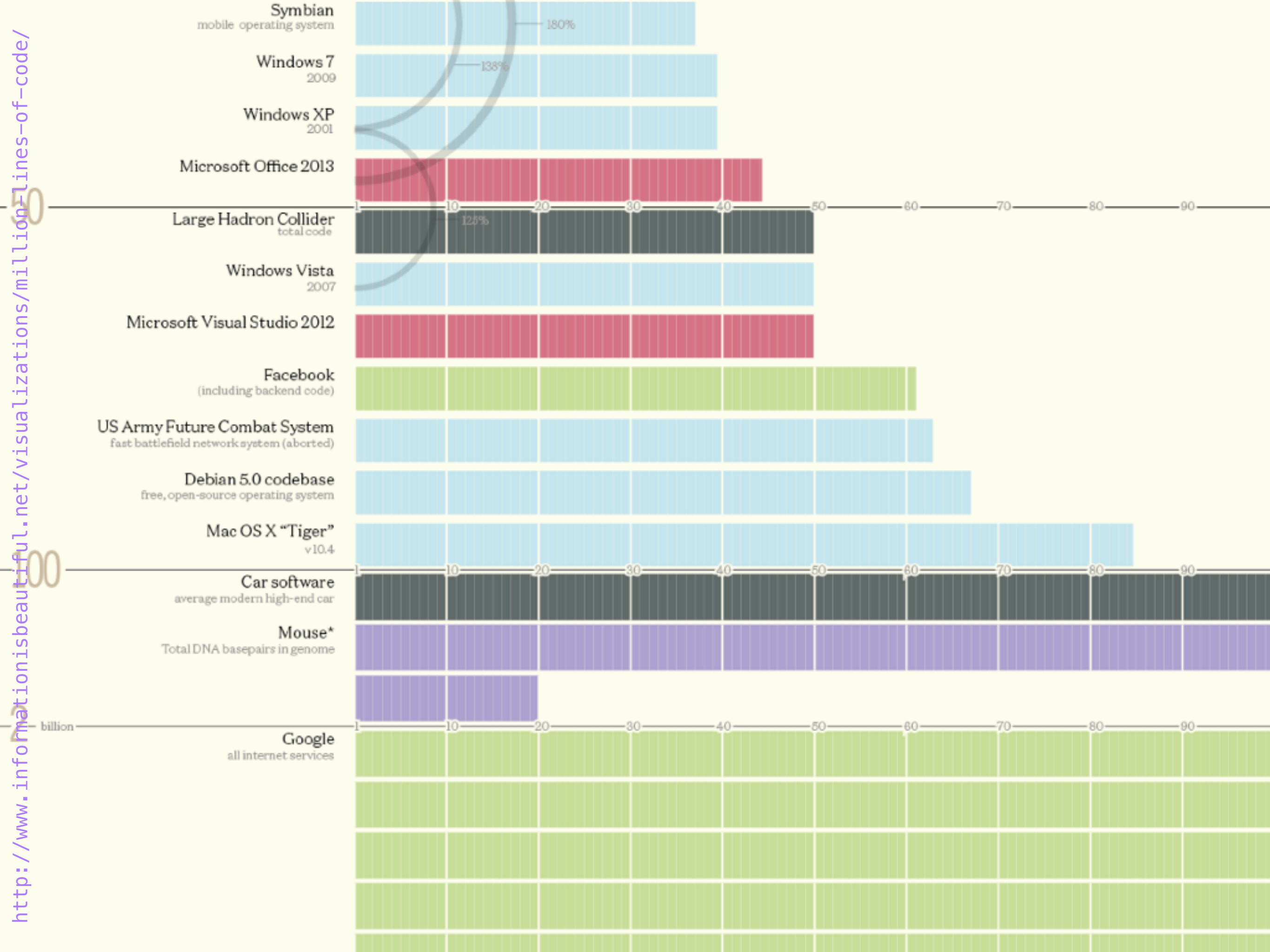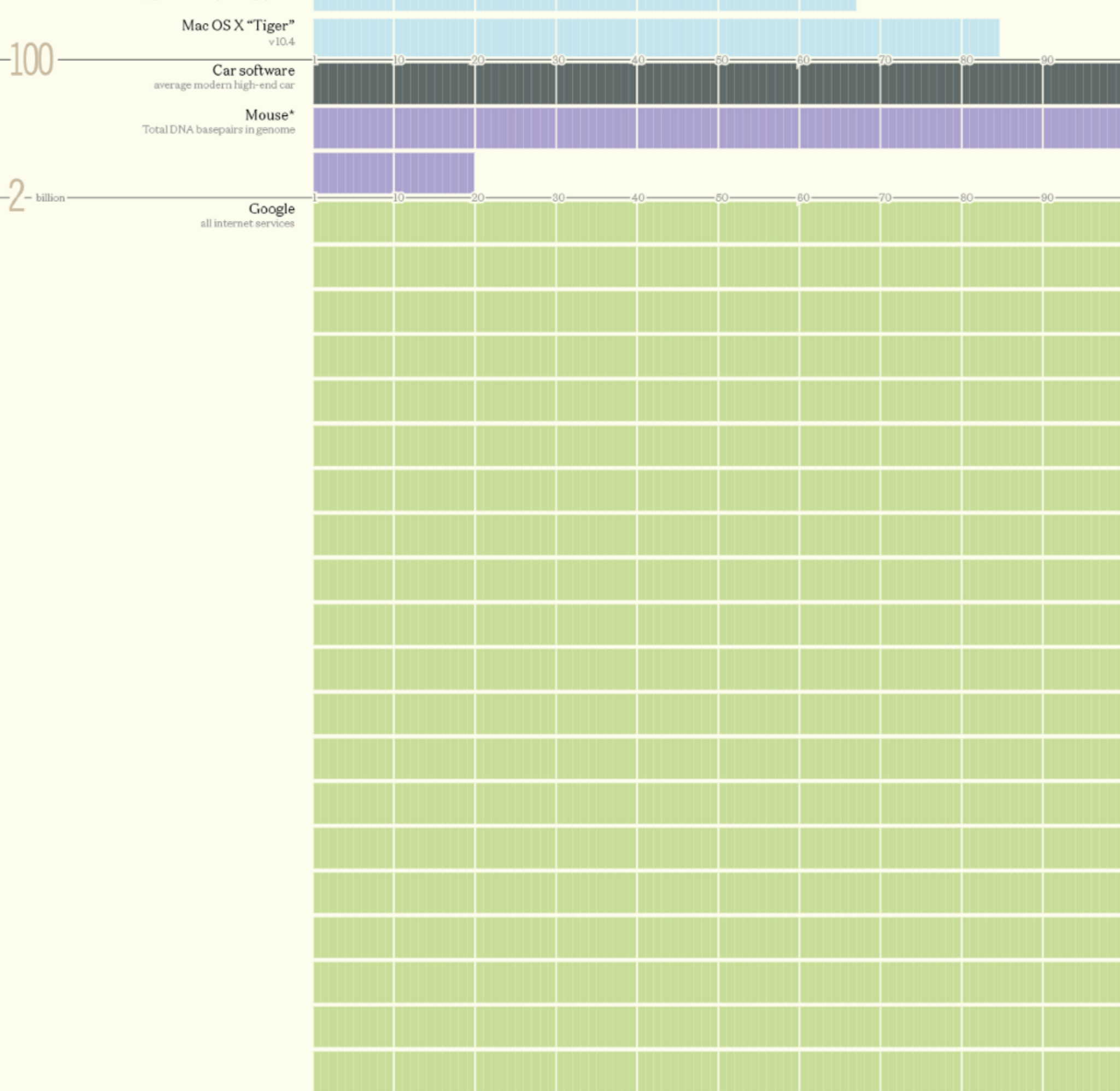hundred thousand

0  1  2  3  4  5  6  7  8  9

simple iPhone game app

Unix v 1.0
1971

OPERATING SYSTEM

Win32/Simile virus

average iPhone app

APP

Pacemaker

Photoshop v. 1.0
1990

Camino
web browser

BROWSER

Quake 3 engine
3D Video game system

GAME

Space Shuttle

MACHINE

a million lines of code

million

1  10  20  30  40  50  60  70  80  90

a million lines of code
18,000 pages of printed text

War And Peace x 14, or
Ulysses x 25, or
The Catcher in The Rye x 63

CryEngine 2
3D video game system

Bacteria
Syphillis (Treponema pallidum)

ORGANISM

Age of Empires online

3780%

CESM Climate Model
National Center for Atmospheric Research

CESM Climate Model
National Center for Atmospheric Research

F-22 Raptor fighter jet

Linux Kernel 2.2.0
core code

Jurassic Park codebase
source: Dennis Nedry

Hubble Space Telescope

Unreal engine 3
3D video game system

Windows 3.1
1992

Large Hadron Collider
(root software)

US military drone
(control software only)

288%

Photoshop C.S. 6
image editing software

Windows NT 3.1
1993

HD DVD Player on XBox
(just the player)

needed to repair HealthCare.gov
apparently

WEB

Mars Curiosity Rover
Martian ground vehicle probe

Linux kernel 2.6.0
2003

166%

Google Chrome
latest

upper
estimate

World of WarCraft
server only

Boeing 787
avionics & online support systems only

Windows NT 3.5

1    10    20    30    40    50    60    70    80    90

**Windows NT 3.5**
1993

**Firefox**
latest version

**Chevy Volt**
electric car

**Intuit Quickbooks**
accounting software

**Windows NT 4.0**
1996

**Android**
mobile device operating system

**Mozilla Core**
core code at heart of all Mozilla's software

**MySQL**
database language

**Boeing 787**
total flight software

**Linux 3.1**
recent version

**Apache Open Office**
open-source office software

**F-35 Fighter jet**
2013

**Microsoft Office 2001**

**Windows 2000**

**Microsoft Office for Mac**
2006

**Symbian**
mobile operating system

**Windows 7**
2009

**Windows XP**
2001

10

288%

252%

1  10  20  30  40  50  60  70  80  90

25

1  10  20  30  40  50  60  70  80  90

180%

138%

Symbian
mobile operating system

180%

Windows 7
2009

138%

Windows XP
2001

Microsoft Office 2013

Large Hadron Collider
total code

125%

Windows Vista
2007

Microsoft Visual Studio 2012

Facebook
(including backend code)

US Army Future Combat System
fast battlefield network system (aborted)

Debian 5.0 codebase
free, open-source operating system

Mac OS X "Tiger"
v10.4

Car software
average modern high-end car

Mouse*
Total DNA basepairs in genome

Google
all internet services

Mac OS X "Tiger"
v10.4

100

Car software
average modern high-end car

Mouse*
Total DNA basepairs in genome

2 — billion

Google
all internet services

1    10    20    30    40    50    60    70    80    90

# Schedule

| W44 | Introduction | V.Zaytsev |
|-----|--------------|-----------|
| W45 | Metaprogramming | J.Vinju |
| W46 | Reverse Engineering | V.Zaytsev |
| W47 | Software Analytics | M.Bruntink |
| W48 | Clone Management | M.Bruntink |
| W49 | Source Code Manipulation | V.Zaytsev |
| W50 | Legacy and Renovation | TBA |
| W51 | Conclusion | V.Zaytsev |

# Schedule

| W44 | Introduction | V.Zaytsev |
| --- | --- | --- |
| W45 | Metaprogramming | J.Vinju |
| W46 | Reverse Engineering | V.Zaytsev |
| W47 | Software Analytics | M.Bruntink |
| W48 | Clone Management | M.Bruntink |
| W49 | Source Code Manipulation | V.Zaytsev |
| W50 | Legacy and Renovation | TBA |
| W51 | Conclusion | V.Zaytsev |

# Schedule



| W44 | Introducti... | V.Zaytsev |
| W45 | Metaprogramm... | J.Vinju |
| W... | ...neering | V.Zaytsev |
| W... | ...lytics | M.Bruntink |
| W48 | Clone Manage... | |
| W49 | Source Code Manipulation | V.Zaytsev |
| W50 | Legacy and Renovat... | |
| W51 | | V.Zaytsev |

**Series 1:** implement a set of metrics

**Series 2:** write a clone detector

Review a paper

Honours Track

# Deadlines & Deliverables

* 2 Nov:  Series 0 (Rascal test)

* 17 Nov: Series 1 = $\frac{1}{3}$ grade

* 1 Dec:  Review = $\frac{1}{3}$ grade

* 15 Dec: Series 2 = $\frac{1}{3}$ grade

# Teachers

Dr. Vadim Zaytsev

Dr. Magiel Bruntink

Prof.Dr. Jurgen Vinju

Davy Landman

SWAT

Jouke Stoel

# Software Types

# Program Types: S

* S-type programs

  * "specifiable"

  * problem formally defined by a spec

  * automated acceptance possible

  * such software does not evolve

M.M.Lehman, *Programs, Life Cycles and Laws of Software Evolution*, IEEE 68(9), 1980.

# Program Types: S



**S-type**

# Program Types: P

* P-type programs

  * "problem-solving"

  * problem models a real-world task

    * imperfectly

  * qualitative acceptance

  * they can evolve continuously

M.M.Lehman, *Programs, Life Cycles and Laws of Software Evolution*, IEEE 68(9), 1980.

# Program Types: P



**P-type**

# Program Types: E

* E-type programs

  * "embedded"

  * solution is a part of the world

  * acceptance is subjective

  * they are inherently evolutionary

M.M.Lehman, *Programs, Life Cycles and Laws of Software Evolution*, IEEE 68(9), 1980.

# Program Types: E

# Lehman's Laws of Software Evolution

# Lehman's Laws (1/8)

* Continuing Change

  * E-system rots unless adapted

  * the process never stops

  * (true for P-systems as well)

M.M.Lehman, *Programs, Life Cycles and Laws of Software Evolution*, IEEE 68(9), 1980.

# Lehman's Laws (2/8)

* Increasing Complexity
  * E-system becomes more complex
  * evolving means complicating
  * (unless we do something)

M.M.Lehman, *Programs, Life Cycles and Laws of Software Evolution*, IEEE 68(9), 1980.

# Lehman's Laws (3/8)

* Self-regulation

  * E-system evolution is SRP

  * obeys certain statistical laws

  * (distribution close to normal)

M.M.Lehman, *Programs, Life Cycles and Laws of Software Evolution*, IEEE 68(9), 1980.

# Lehman's Laws (4/8)

* **Conservation of Organisational Stability**
  * E-system dev activity is invariant
  * throughout its lifetime
  * (does not depend on resources)

# Lehman's Laws (5/8)

* **Conservation of Familiarity**
  * E-system changes per release
    * invariant
  * throughout its lifetime
  * (too little: bored;
    too much: overwhelmed)

M.M.Lehman, *Programs, Life Cycles and Laws of Software Evolution*, IEEE 68(9), 1980.

# Lehman's Laws (6/8)

* Continuing Growth

  * E-system must add features over time

  * to keep users satisfied

  * (expectations creep)

M.M.Lehman, J.F.Ramil, P.D.Wernick, D.E.Perry, W.M.Turski, *Metrics and Laws of Software Evolution — The Nineties View*, METRICS, 1997.

# Lehman's Laws (7/8)

* **Declining Quality**
  * E-system perceived quality declines
  * internal as well as external
  * (unless constantly maintained)

M.M.Lehman, J.F.Ramil, P.D.Wernick, D.E.Perry, W.M.Turski, *Metrics and Laws of Software Evolution — The Nineties View*, METRICS, 1997.

# Lehman's Laws (8/8)

* Feedback System

  * E-system evolution is a

    * feedback system

  * multi-level

  * multi-loop

  * multi-agent

M.M.Lehman, J.F.Ramil, P.D.Wernick, D.E.Perry, W.M.Turski, *Metrics and Laws of Software Evolution — The Nineties View*, METRICS, 1997.

# Lehman's Laws

* Continuing Change

* Increasing Complexity

* Self-regulation

* Conservation of Organisational Stability

* Conservation of Familiarity

* Continuing Growth

* Declining Quality

* Feedback System

M.M.Lehman, J.F.Ramil, P.D.Wernick, D.E.Perry, W.M.Turski, *Metrics and Laws of Software Evolution – The Nineties View*, METRICS, 1997.

# Maintenance Types

# Maintenance

* Modification of a software product after delivery to correct faults, to improve performance or other attributes, or to adapt the product to a modified environment

IEEE 1219, 1993

# Maintenance phases

* Introductory

  * user support!

* Growth

  * correcting faults!

* Maturity

  * enhancements!

* Decline

  * technology replacement!

Hans van Vliet, *Software Engineering: Principles and Practice*. Jon Wiley & Sons, 2009.

# Types of maintenance

* Corrective

* Adaptive

* Perfective

* Preventive

B.P.Lientz, E.B.Swanson, *Software Maintenance Management, A Study of the Maintenance of Computer Application Software in 487 Data Processing Organizations*, 1980.

# Types of maintenance

* Corrective

* Adaptive

* Perfective

* Preventive

B.P.Lientz, E.B.Swanson, *Software Maintenance Management, A Study of the Maintenance of Computer Application Software in 487 Data Processing Organizations*, 1980.

# Types of maintenance

* Corrective

* Adaptive

* Perfective

  * user enhancement

  * efficiency

  * other

* Preventive



B.P.Lientz, E.B.Swanson, *Software Maintenance Management, A Study of the Maintenance of Computer Application Software in 487 Data Processing Organizations*, 1980.

# Top 5 problems

∗ Quality of documentation

∗ User demand for enhancements

∗ Competing demands for maintainers' time

∗ Meeting scheduled commitments

∗ Turnover in user organisations

S.L.Pfleeger, *Software Engineering: Theory and Practice,* Prentice Hall, 1998.

# Is it hopeless?

* Higher quality

  * less (c) maintenance

* Anticipating changes

  * less (a&p) maintenance

* Better tuning to user needs

  * less (p) maintenance

* Less code

  * less (*) maintenance

# Roadmap

∗ Metaprogramming

∗ Reverse engineering

∗ Software analytics

∗ Clone management

∗ Source code manipulation

∗ Legacy

# State of the Art

# Detecting Complex Changes During Metamodel Evolution

* Metamodel evolves:





* Follow user actions

* Detect complex patterns

* Enrich evolution trace



## Detecting Complex Changes During Metamodel Evolution

Djamel Eddine Khelladi[1]([✉]), Regina Hebig[1], Reda Bendraou[1],
Jacques Robin[1], and Marie-Pierre Gervais[1,2]

[1] Sorbonne Universités, UPMC Univ Paris 06, UMR 7606, F-75005 Paris, France
djamel.khelladi@lip6.fr
[2] Université Paris Ouest Nanterre La Defense, F-92001 Nanterre, France

**Abstract.** Evolution of metamodels can be represented at the finest grain by the trace of atomic changes: add, delete, and update elements. For many applications, like automatic correction of models when the metamodel evolves, a higher grained trace must be inferred, composed of complex changes, each one aggregating several atomic changes. Complex change detection is a challenging task since multiple sequences of atomic changes may define a single user intention and complex changes may overlap over the atomic change trace. In this paper, we propose a detection engine of complex changes that simultaneously addresses these two challenges of variability and overlap. We introduce three ranking heuristics to help users to decide which overlapping complex changes are likely to be correct. We describe an evaluation of our approach that allow reaching full recall. The precision is improved by our heuristics from 63% and 71% up to 91% and 100% in some cases.

**Keywords:** Metamodel · Evolution · Complex change · Detection

### 1 Introduction

In the process of building a domain-specific modeling language (DSML) multiple versions are developed, tried out, and adapted until a stable version is reached. As by one of our industrial partners in the automotive domain, such intermediate versions of the DSML are used in product development, where often further needs are identified. A challenge hereby is that each time the metamodel of the DSML is changed to a next version, already developed models need to be co-evolved too. This is not only the case for DSMLs, but also for more generic metamodels, e.g. the UML officially evolved in the past every two to three years.

To cope with this evolution of metamodels, mechanisms are developed to co-evolve artifacts, such as models and transformations that may become invalid. A challenging task herein is to detect all the changes that lead a metamodel from a version $n$ to a version $n+1$, called Evolution Trace (ET). Automatically detecting it, not only helps developers to automatically keep track of the metamodels' evolution, but also to trigger and/or to apply automatic actions based on these changes. For instance, models and transformations that are defined based on the metamodel are automatically co-evolved i.e. corrected based on the detected

http://bibtex.github.io/CAISE-2015-KhelladiHBRG.html
http://dx.doi.org/10.1007/978-3-319-19069-3_17

# Automated Unit Test Generation for Evolving Software

* Software evolves.

* How about test cases?

* Functionality changes?

   * (regression testing)

* Tests that used to work

generated
test suite

old
version .java

EVOSUITER

new
version .java

Fitness
Function
- coverage
- state distance
- control-flow distance

old
version .java

new
version .java

Automated Unit Test Generation for Evolving Software

Sina Shamshiri
Department of Computer Science, University of Sheffield
Regent Court, 211 Portobello, Sheffield, UK, S1 4DP
sina.shamshiri@sheffield.ac.uk

**ABSTRACT**
As developers make changes to software programs, they want to ensure that the originally intended functionality of the software has not been affected. As a result, developers write tests and execute them after making changes. However, high quality tests are needed that can reveal unintended bugs, and not all developers have access to such tests. Moreover, since tests are written without the knowledge of future changes, sometimes new tests are needed to exercise such changes. While this problem has been well studied in the literature, the current approaches for automatically generating such tests either only attempt to reach the change and do not aim to propagate the infected state to the output, or may suffer from scalability issues, especially when a large sequence of calls is required for propagation. We propose a search-based approach that aims to automatically generate tests which can reveal functionality changes, given two versions of a program (e.g., pre-change and post-change). Developers can then use these tests to identify unintended functionality changes (i.e., bugs). Initial evaluation results show that our approach can be effective on detecting such changes, but there remain challenges in scaling up test generation and making the tests useful to developers, both of which we aim to overcome.

**Categories and Subject Descriptors**
D.2.5 [**Software Engineering**]: Testing and Debugging—*Testing Tools*; I.2.8 [**Artificial Intelligence**]: Problem Solving, Control Methods, and Search

**Keywords**
Automated Unit Test Generation, Genetic Algorithms, Search-Based Testing, Regression Testing

**1. INTRODUCTION**
Developers evolve software programs by introducing many changes throughout the life-cycle of the software. These changes often range from small refactorings to the addition of large new features. However, some of these changes may affect the originally intended functionality of the software, by introducing unintended bugs – also known as *regression faults*. To avoid regressions in the functionality, engineers write tests as they develop the software, and after making changes developers execute these tests to increase their confidence that the intended functionality of the software is intact. This practice is also referred to as *regression testing* and is commonly used in the industry.

While regression testing can help with early detection of regression faults, developers face several challenges when applying the technique. As the number of tests grows, execution of all tests after every single change can become expensive and impractical. This problem has been well studied in the literature [18] and many techniques such as test selection, prioritization and minimization have been proposed.

The challenges however are not limited to the growing cost of regression testing. Even if all tests are executed, three main problems remain: 1) an existing set of tests is required, 2) the tests are often written without foreseeing future changes, and 3) the effectiveness of the tests in finding regression faults depends on the quality of the written tests. According to the PIE model [15], to reveal a fault, a test has to first execute the fault, infect the state and finally propagate it to the output. While several techniques exist for augmenting existing test suites (e.g., [10,17]) and generating regression tests (e.g., [2,9,13,14]), the techniques mainly focus on reaching the fault, yet the number of paths to propagate the infected state to the output can explode, which may impose a limit on the scalability of the approach [3].

To address the previous shortcomings, we propose a technique for generating a regression test suite (i.e. a set of unit-tests which contain a sequence of calls executing the class under test) without depending on existing tests. Our approach takes two versions of a class under test, and uses a search-based algorithm [8] with the objective of reaching and propagating the changes between the two versions of the program. We have implemented our approach named EVOSUITE on top of the EvoSuite [5] test generation tool, and our early evaluation of the technique [11] showed encouraging results on examples with propagation issues (i.e. where covering the change alone does not propagate the changed state to the output). Further attempts to evaluate the effectiveness of our approach on detecting real regression faults revealed several challenges. As a part the remaining course of this research we aim to solve these challenges, in addition to evaluating our approach against the state-of-the-art.

1038

# Detection of Software Evolution Phases based on Development Activities

* Software evol. history:

  * commits — fine-grained

  * releases — coarse

* Something in between?

* 8 kinds of phases:

  * changes: important/not

  * dev: rapid/slow

  * change types: different/same



A: Important + Rapid + Different    K: Less-Imp. + Rapid + Different
B: Important + Rapid + Similar      L: Less-Imp. + Rapid + Similar
C: Important + Slow + Different      M: Less-Imp. + Slow + Different
D: Important + Slow + Similar        N: Less-Imp. + Slow + Similar

% of release duration

http://bibtex.github.io/ICPC-2015-BenomarASPS.html
http://dx.doi.org/10.1109/ICPC.2015.11

# *Evolution of Software Development Strategies*

* Expert/novice devs

* Look at students

  * first year

  * final year

# *A Study on the Role of Software Architecture in the Evolution and Quality of Software*

* Impact of architecture

  on evolution?

* Problem: documentation

* Reverse engineer!

  * heuristics, IR, DM, …

  * e.g. modularisation

* Architectural bad smells

  * x-module co-change

  * 20% defects, 2× time2fix



http://bibtex.github.io/MSR-2015-KouroshfarMBXMC.html
http://dx.doi.org/10.1109/MSR.2015.30

# *Modelling the Evolution of Development Topics using Dynamic Topic Models*

* Tasks evolve with sw

* Can be grouped by topic

  * Strength evolution

  * Content evolution

  * (never together)

* Use unstructured repos

* Visualise!



http://bibtex.github.io/SANER-2015-HuSLL.html
http://dx.doi.org/10.1109/SANER.2015.7081810

# *Modelling the Evolution of Development Topics using Dynamic Topic Models*

* Tasks evolve with sw

* Can be grouped by topic

# *Mining Software Contracts for Software Evolution*

* Version control systems

    * git, svn, cvs

* Program contract

    * pre- & postcond, inv

    * "requires", "ensures"…

* Bugfixs & contracts coevolve

* Some things easier to RE

    * from contracts

* http://github.com/ybank/inv-research

http://bibtex.github.io/ICSME-2014-YanMG.html
http://dx.doi.org/10.1109/ICSME.2014.76

# *Visualising the Evolution of Systems and Their Library Dependencies*

# *Visualising the Evolution of Systems and Their Library Dependencies*

# *Visualising the Evolution of Systems and Their Library Dependencies*



**Library Version Usage**

- 1.00
- 0.75
- 0.50
- 0.25

**Library Evolution Types**

- \+ adopter
- · dropped
- o idler
- ◇ updater

# Conclusion

* Software evolves

* Software evolution obeys certain laws

* Software rots in time (quality, complexity…)

* 70% of software engineers do maintenance

* Many software systems are legacy

* Forward, reverse and re-engineering

* Actively researched field

* Learn to build tools