

Models for hierarchical inheritance structures in object-oriented programming languages

Mariacristina Romano

UNIVERSITÀ DEGLI STUDI DI MILANO
FACOLTÀ DI SCIENZE E TECNOLOGIE
Corso di laurea magistrale in Fisica

15 aprile 2015

Complex Systems and Computer science

Introduction

Complex Systems and Computer science

Object-oriented paradigm

Introduction

Complex Systems and Computer science

Object-oriented paradigm

Inheritance

Introduction

Complex Systems and Computer science

Object-oriented paradigm

Inheritance

Code reuse

Introduction

Complex Systems and Computer science

Object-oriented paradigm

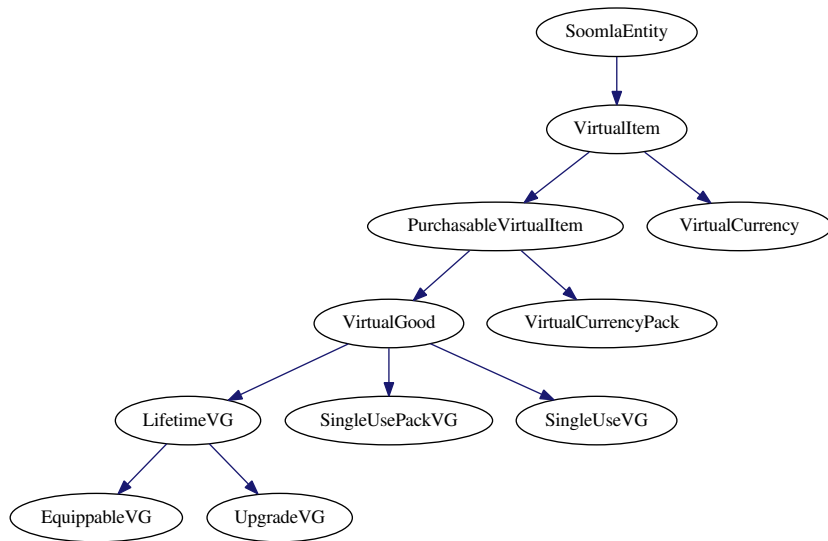
Inheritance

Code reuse

Hierarchies

Example of inheritance hierarchy

Data Analysis - Project Soomla Cocos2dx



Noisy complex system dataset

Data Analysis - dataset

Packages have been downloaded from [GitHub](#), the actual largest code host on the web.

To give a **complete overview** of inheritance hierarchies, three different programming languages have been analyzed.

The **dataset** contains:

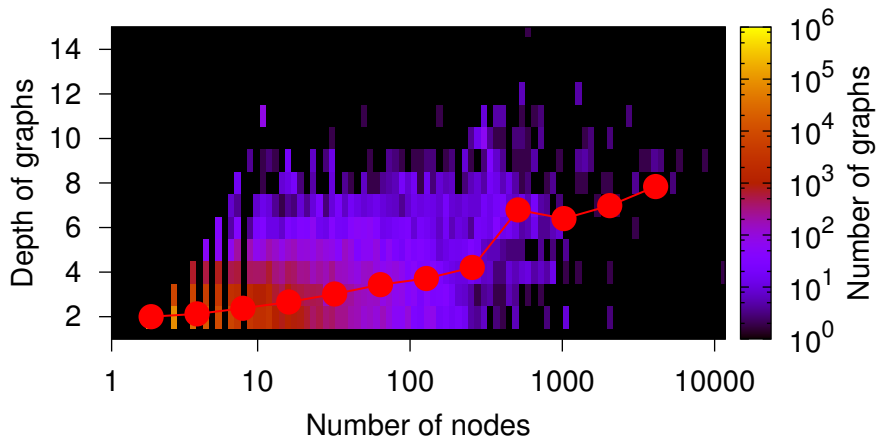
- 17333 C++ projects (3233447 hierarchies)
- 25318 Java projects (3504681 hierarchies)
- 20010 Python projects (2491603 hierarchies)

Almost 10 millions of inheritance hierarchies!

Depth VS Size is logarithmic

Data Analysis

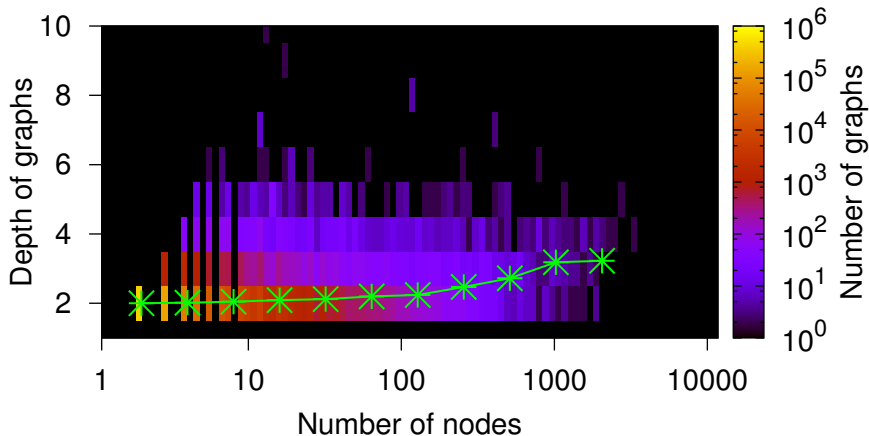
C++



Depth VS Size is logarithmic

Data Analysis

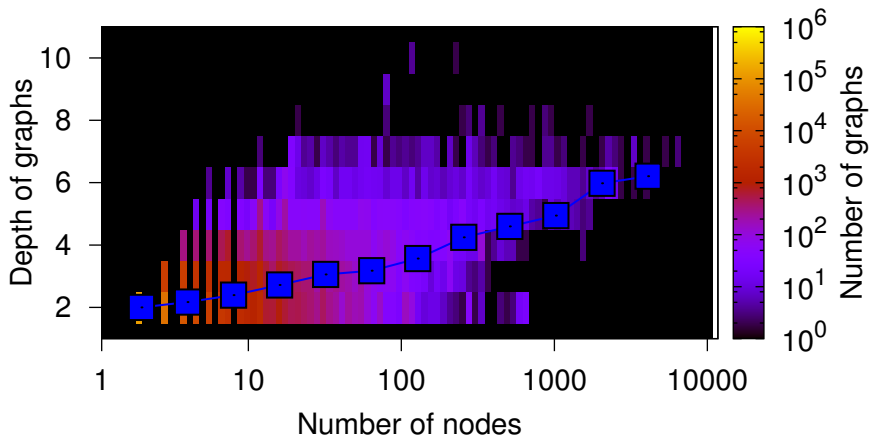
Java



Depth VS Size is logarithmic

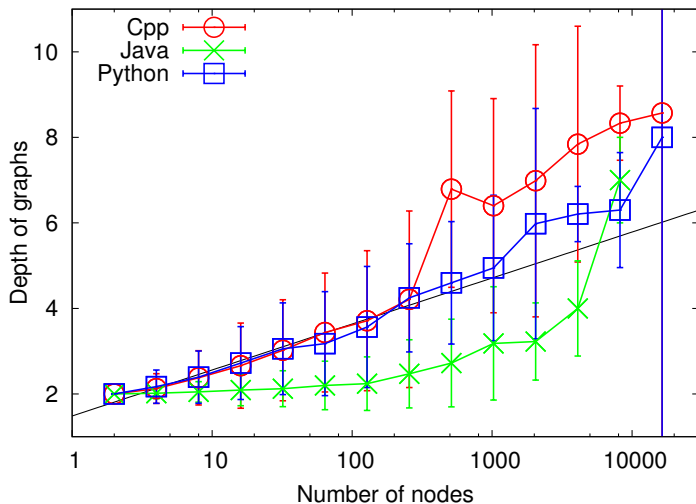
Data Analysis

Python



Depth VS Size is logarithmic

Data Analysis - Comparison among languages



Sharing Tree model

(A microscopic model)

How to build the structure

Sharing Tree model (microscopic model)



create
stars
polygons

draw
freehand
lines

create
text
objects

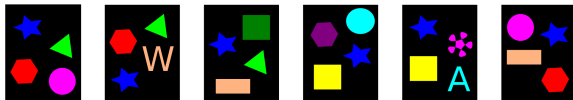
fill
bounded
areas

pick
colors

erase
existing
paths

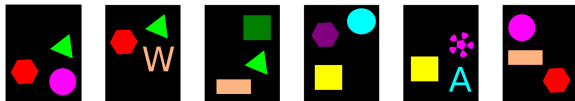
How to build the structure

Sharing Tree model (microscopic model)



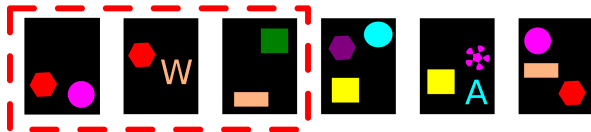
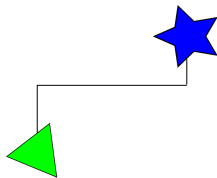
How to build the structure

Sharing Tree model (microscopic model)



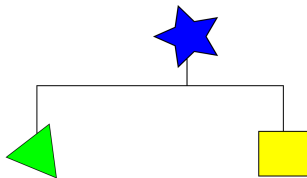
How to build the structure

Sharing Tree model (microscopic model)



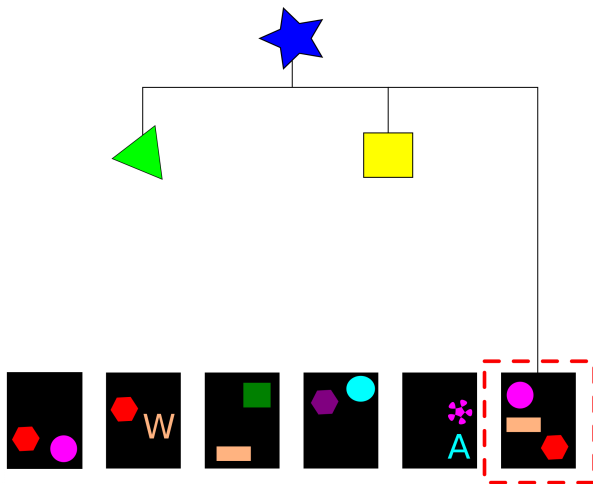
How to build the structure

Sharing Tree model (microscopic model)



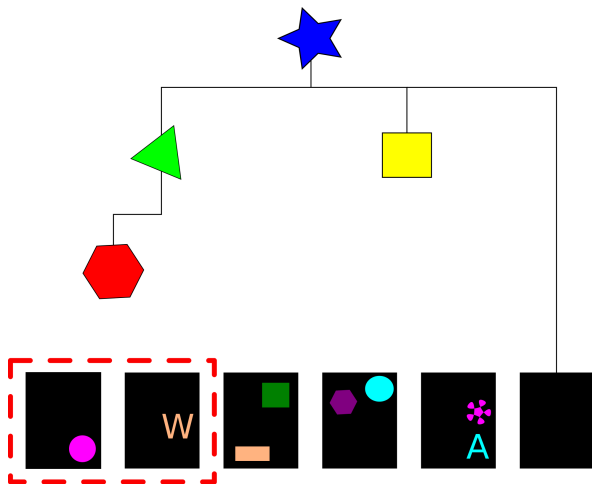
How to build the structure

Sharing Tree model (microscopic model)



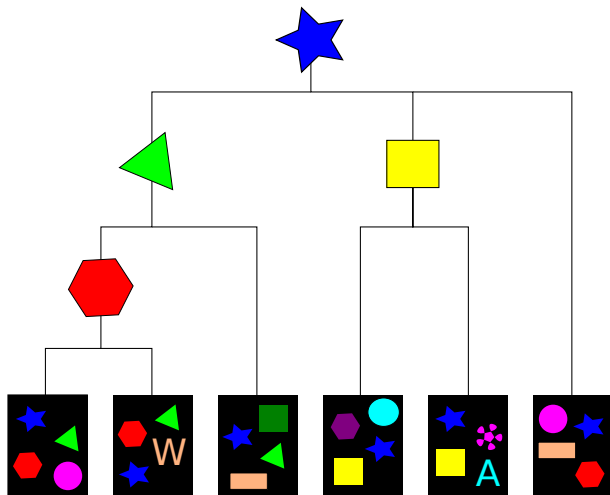
How to build the structure

Sharing Tree model (microscopic model)



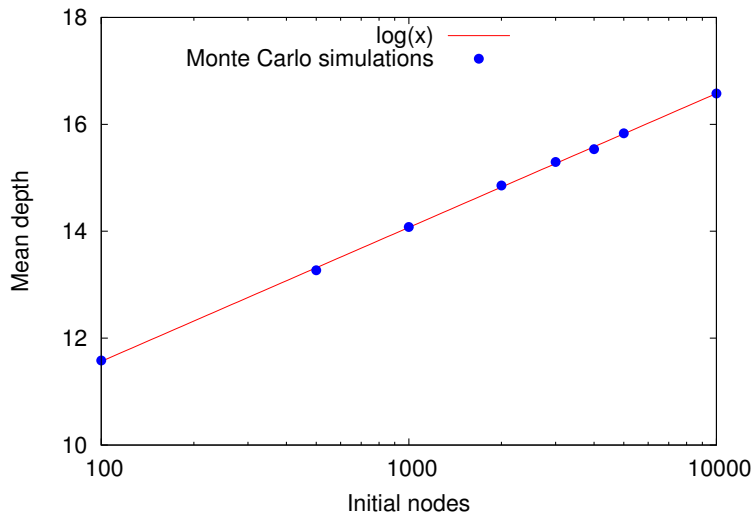
How to build the structure

Sharing Tree model (microscopic model)



Depth VS Size is logarithmic

Sharing Tree model (microscopic model)



Minimal Effort model

(A mean field model)

The Effort to build a hierarchy

Minimal Effort model (mean field model)

$$E = \sum_{\sigma}^{\mathcal{N}} \text{cost}(\sigma)$$

You need n classes to perform a task

Minimal Effort model (mean field model)



create
stars
polygons

draw
freehand
lines

create
text
objects

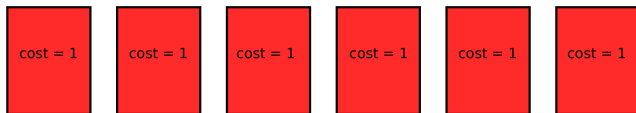
fill
bounded
areas

pick
colors

erase
existing
paths

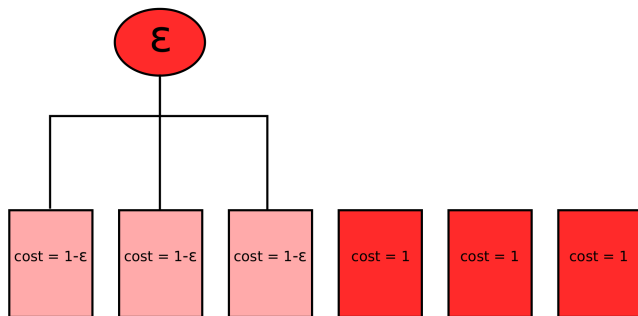
The cost of each class

Minimal Effort model (mean field model)



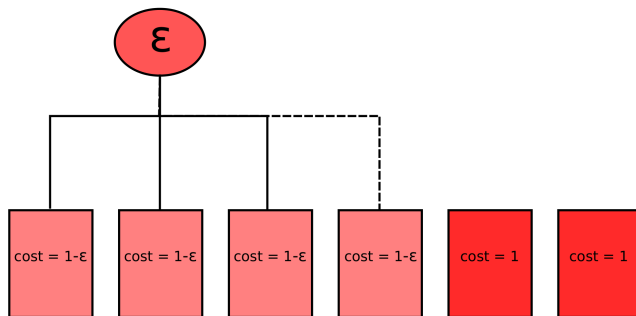
Reuse

Minimal Effort model (mean field model)



Competition

Minimal Effort model (mean field model)



Competition

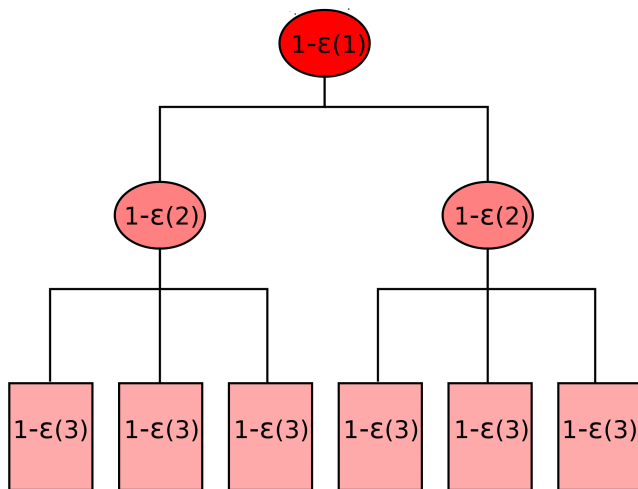
Minimal Effort model (mean field model)

$$E = \sum_{\sigma}^{\mathcal{N}} \text{cost}(\sigma)$$

$$E = \sum_{\sigma}^{\mathcal{N}} [1 - \varepsilon(\mathbf{m})]$$

The effort of “writing” a tree

Minimal Effort model (mean field model)



How much of the code is shareable?

Minimal Effort model (mean field model)

The probability to find a selected symbol in a sequence of k extractions is

$$p = 1 - \left(1 - \frac{1}{S}\right)^k$$

How much of the code is shareable?

Minimal Effort model (mean field model)

The probability to find a selected symbol in a sequence of k extractions is

$$p = 1 - \left(1 - \frac{1}{S}\right)^k$$

$$S \rightarrow \infty \quad k \rightarrow \infty \quad \beta \equiv \frac{k}{S} \quad e^{-\beta} = \lim_{S \rightarrow +\infty} \left(1 - \frac{1}{S}\right)^{\beta S}$$

How much of the code is shareable?

Minimal Effort model (mean field model)

The probability to find a selected symbol in a sequence of k extractions is

$$p = 1 - \left(1 - \frac{1}{S}\right)^k$$

The probability to find the symbol in m independent sets

$$p = 1 - e^{-\beta} \quad \rightarrow \quad \rho = \left(1 - e^{-\beta}\right)^m$$

How much of the code is shareable?

Minimal Effort model (mean field model)

The probability to find a selected symbol in a sequence of k extractions is

$$p = 1 - \left(1 - \frac{1}{S}\right)^k$$

The probability to find the symbol in m independent sets

$$p = 1 - e^{-\beta} \quad \rightarrow \quad \rho = \left(1 - e^{-\beta}\right)^m$$

The shareable code is

$$\varepsilon(m) = \frac{S}{k} \left(1 - e^{-\beta}\right)^m = \frac{1}{\beta} \left(1 - e^{-\beta}\right)^m \equiv e^{-\alpha m}$$

The shared code

Minimal Effort model (mean field model)

$$E = \sum_{\sigma}^{\mathcal{N}} \text{cost}(\sigma)$$

$$E = \sum_{\sigma}^{\mathcal{N}} [1 - \varepsilon(\mathbf{m})]$$

$$E = \sum_{\sigma}^{\mathcal{N}} [1 - e^{-\alpha \mathbf{m}}]$$

Mean field approach

Minimal Effort model (mean field model)

The number of nodes at each level

$$\{n(l)\}_{l=1}^L = \{n(1), n(2), \dots, n(L) \equiv 1\}$$

The mean number of brothers is

$$m(l) = \frac{n(l)}{n(l+1)}$$

The effort as a sum over levels

$$E[L, \{n(l)\}] = \sum_{l=1}^{L-1} \left[1 - \varepsilon \left(\frac{n(l)}{n(l+1)} \right) \right] n(l)$$

E as a function of the structure

Minimal Effort model (mean field model)

$$E = \sum_{\sigma}^{\mathcal{N}} \text{cost}(\sigma)$$

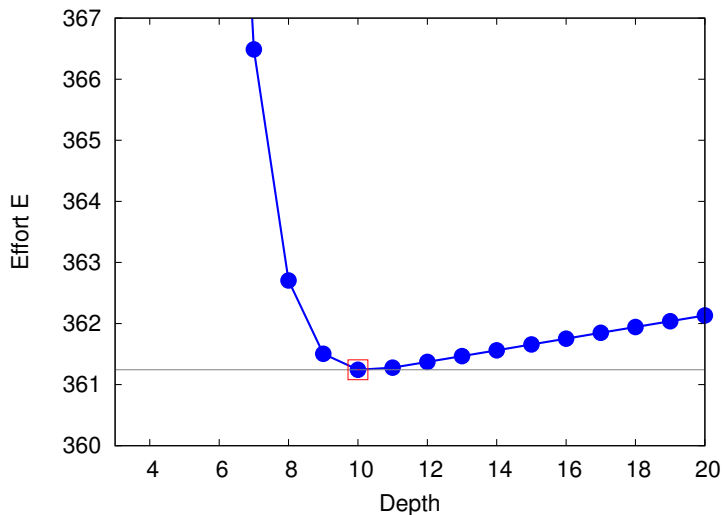
$$E = \sum_{\sigma}^{\mathcal{N}} [1 - \varepsilon(\mathbf{m})]$$

$$E = \sum_{\sigma}^{\mathcal{N}} [1 - e^{-\alpha \mathbf{m}}]$$

$$E[L, \{\mathbf{n}(l)\}] = \sum_{l=1}^{L-1} \left(1 - e^{-\alpha \frac{\mathbf{n}(l)}{\mathbf{n}(l+1)}} \right) \mathbf{n}(l)$$

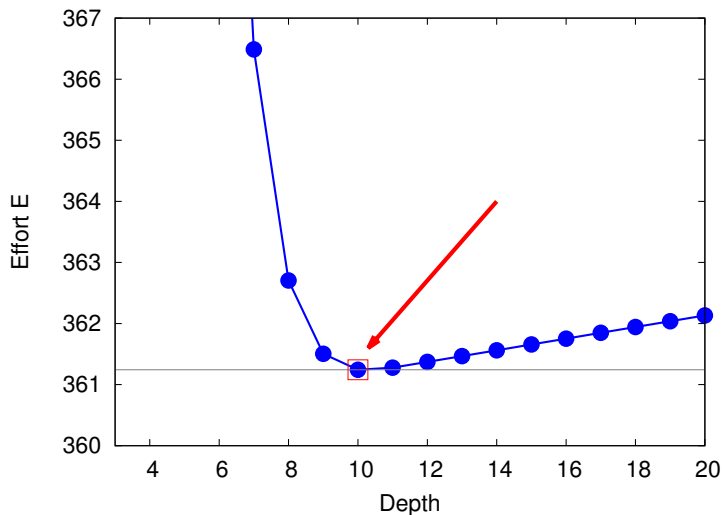
The functional E

Minimal Effort model (mean field model)



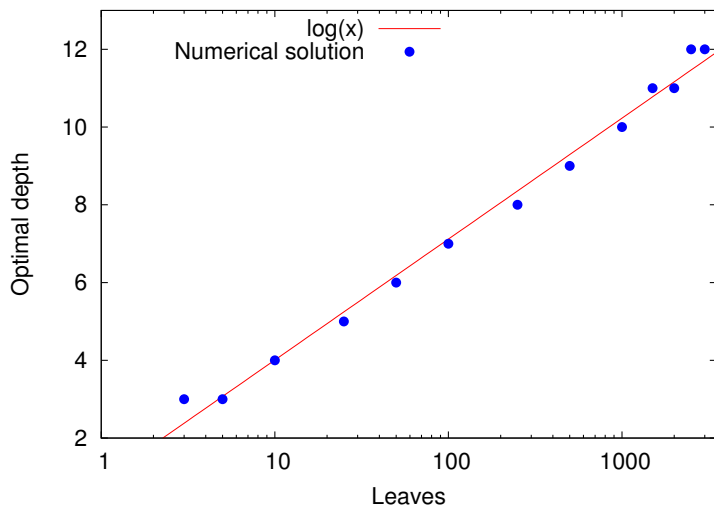
The functional E

Minimal Effort model (mean field model)



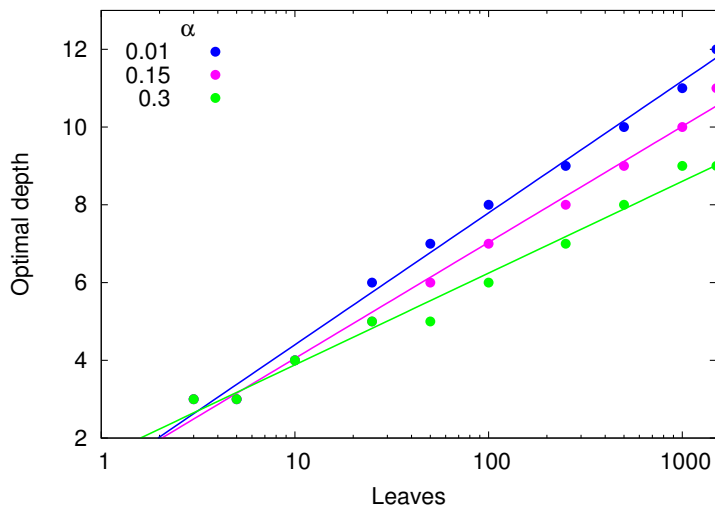
Depth VS Size is logarithmic

Minimal Effort model (mean field model)



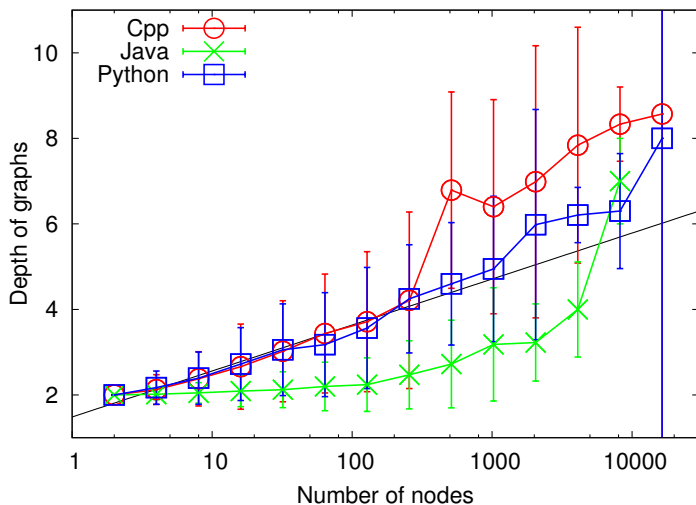
Shareability of the code

Minimal Effort model (mean field model)



Depth VS Size is logarithmic

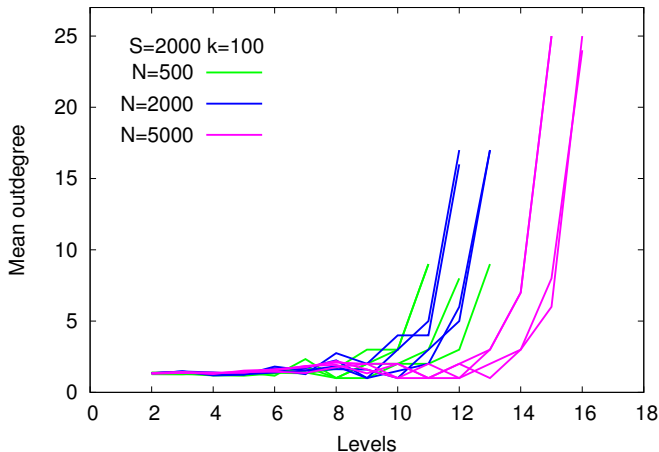
Data Analysis - Comparison among languages



Hierarchies structures

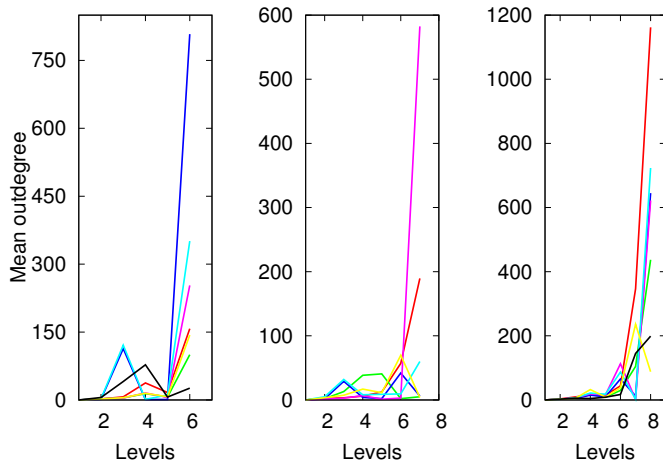
Mean outdegree grows close to the root

Sharing Tree model



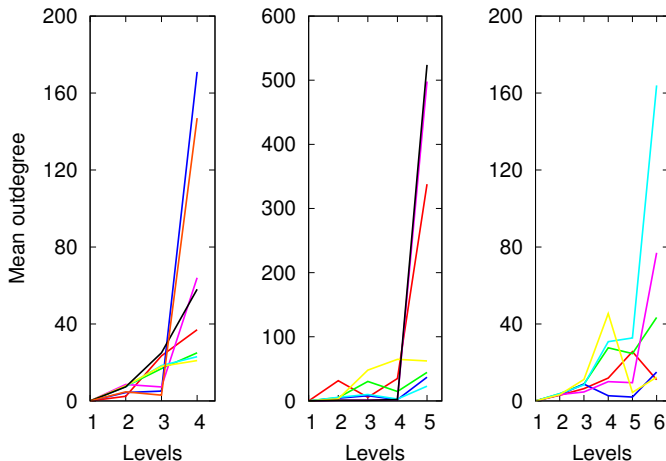
Mean outdegree grows close to the root

C++



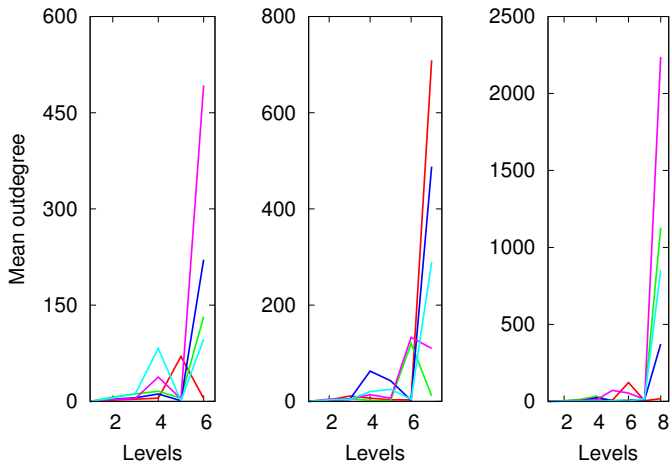
Mean outdegree grows close to the root

Java



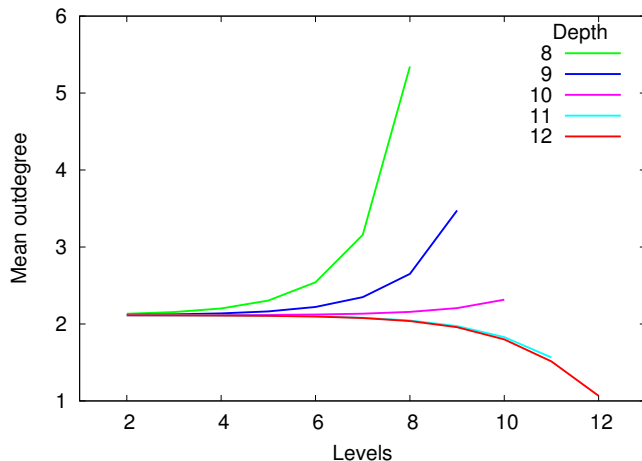
Mean outdegree grows close to the root

Python



Is shallow better?

Minimal Effort model



Conclusions

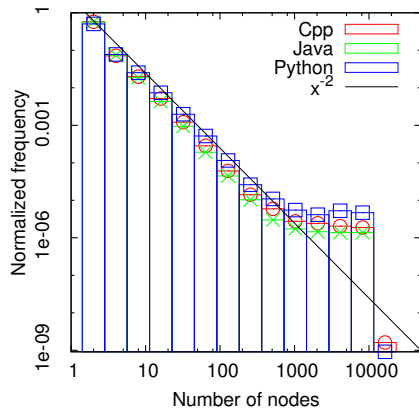
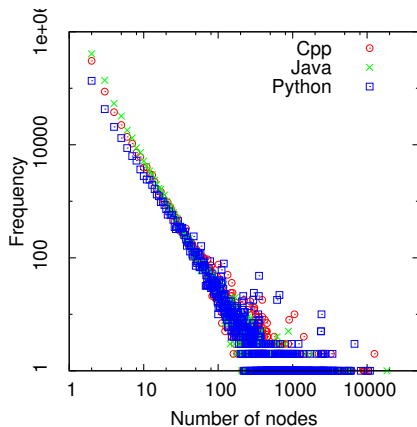
Conclusions

- Different OO programming languages show common behaviors (in sizes distribution, outdegree distribution, depth VS size, ...)
- The two different models (microscopic and mean field) are compatible
- Hierarchies arise from a mechanism of competition between the sake of the reuse and the difficulty of the abstraction
- We have an interpretation about the shallow hierarchies in Java
- Both models predict the growth of the mean outdegree close to the root
- We argue that depths are sub-optimal

Extras

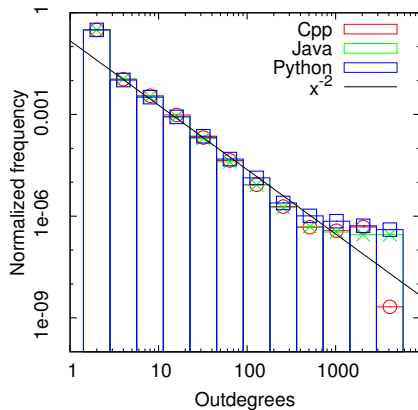
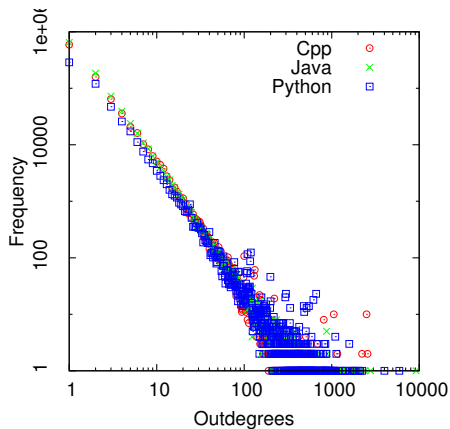
Sizes distribution

Extra



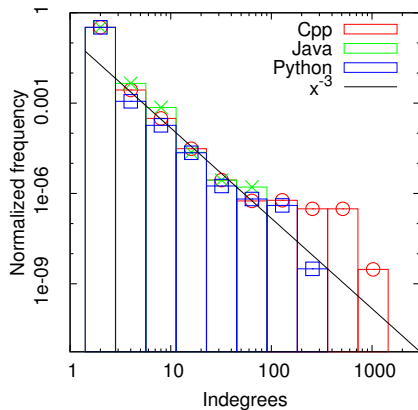
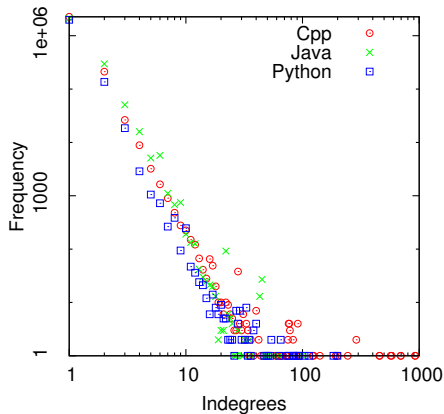
Outdegrees distribution

Extra



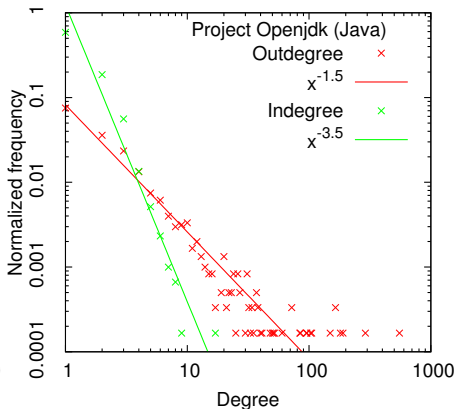
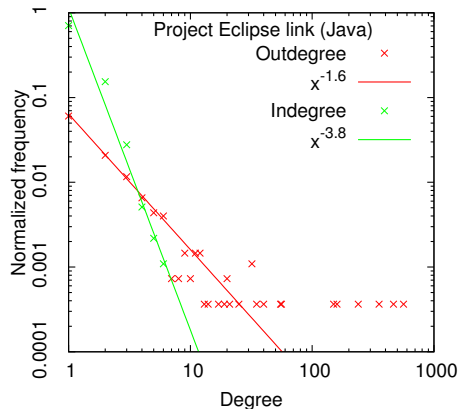
Indegrees distribution

Extra



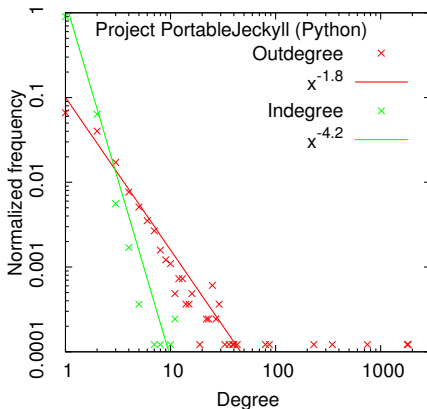
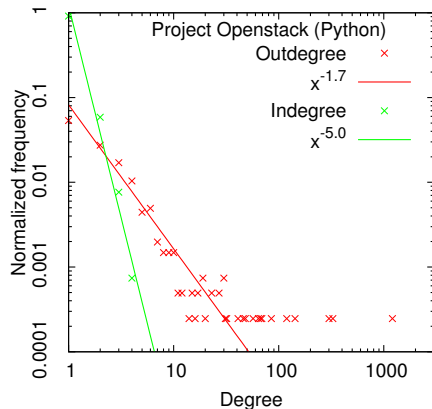
Tree Approximation - Java

Extra



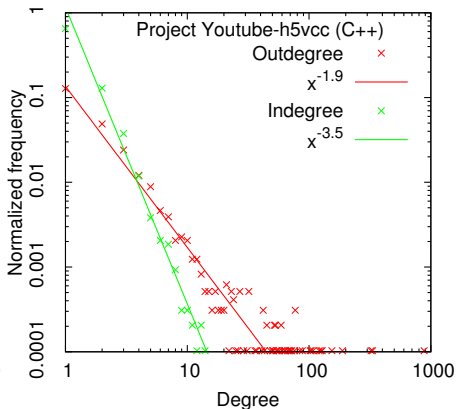
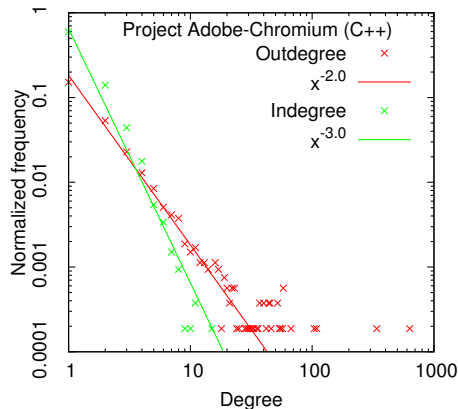
Tree Approximation - Python

Extra



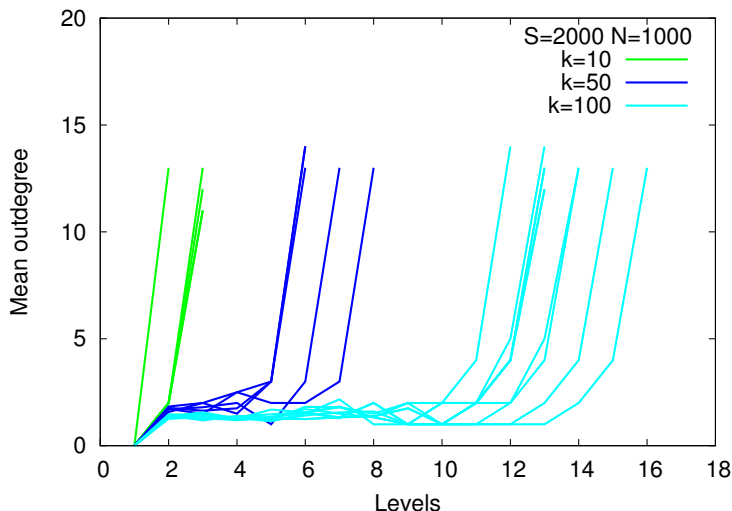
Tree Approximation - C++

Extra



Abstractability in Sharing Tree model

Extra



Most common symbol in Sharing Tree model 1/3

Extra

Probability to find a selected symbol in one sequence is

$$p = \frac{\binom{\mathcal{S}-1}{k-1}}{\binom{\mathcal{S}}{k}} = \frac{k!(\mathcal{S}-1)!}{\mathcal{S}!(k-1)!}$$

Probability that w classes contain a selected symbol with

$$Pr(w) = \binom{\mathcal{N}}{w} p^w (1-p)^{\mathcal{N}-w}$$

Each symbol $s \in \mathcal{S}$ appears in w_s classes. The set $\{w_s\}_{s=1}^{\mathcal{S}}$ contains \mathcal{S} IIDRV. The most common symbol is the one that appears ω times

$$\omega = \max \{w_1, \dots, w_{\mathcal{S}}\}$$

Most common symbol in Sharing Tree model 2/3

Extra

Consider the cumulative distribution function of ω

$$F_{\omega}(y) = Pr(\omega \leq y)$$

Since ω is the maximal occurrence and w_s are independent

$$\begin{aligned} Pr(\omega \leq y) &= Pr(w_1 \leq y, w_2 \leq y, \dots, w_S \leq y) \\ &= Pr(w_1 \leq y)Pr(w_2 \leq y) \dots Pr(w_S \leq y) \end{aligned}$$

and since all w_s have the same cumulative mass function

$$F_{\omega}(y) = F_w^S(y)$$

The probability distribution of ω

$$\begin{aligned} Pr(y = \omega) &= Pr(\omega \leq y) - Pr(\omega \leq y - 1) \\ &= F_w^S(y) - F_w^S(y - 1) \end{aligned}$$

Most common symbol in Sharing Tree model 3/3

Extra

Remembering that w_s are binomial random variables, the occurrence ω of the most common symbol is therefore distributed as

$$\Psi(\omega) = \left(\sum_{i=0}^{\omega} \text{Bin}(\mathcal{N}, i) \right)^S - \left(\sum_{i=0}^{\omega-1} \text{Bin}(\mathcal{N}, i) \right)^S$$

Making explicit the formula of the Binomial distribution

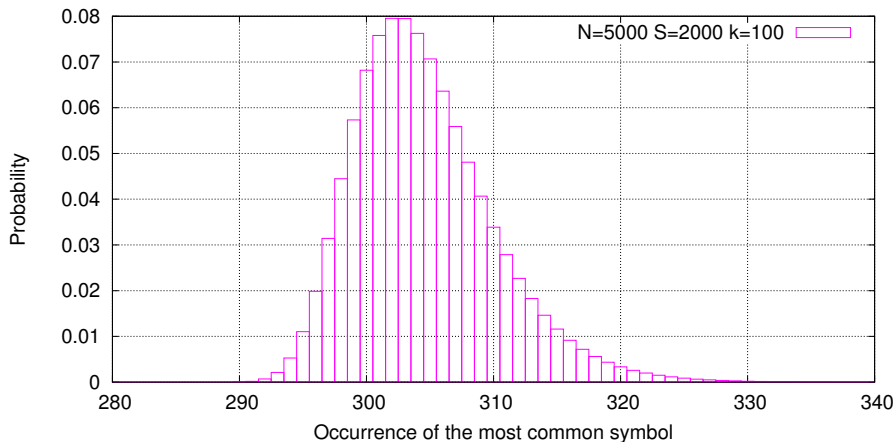
$$\Psi(\omega) = \left(\sum_{i=0}^{\omega} \binom{\mathcal{N}}{i} p^i (1-p)^{\mathcal{N}-i} \right)^S - \left(\sum_{i=0}^{\omega-1} \binom{\mathcal{N}}{i} p^i (1-p)^{\mathcal{N}-i} \right)^S$$

The mean of the distribution is

$$\langle \omega \rangle = \mathcal{N} - \sum_{j=0}^{\mathcal{N}-1} \left(\sum_{k=0}^j \binom{\mathcal{N}}{k} p^k (1-p)^{\mathcal{N}-k} \right)^S$$

Most common symbol distribution

Extra



Sharing Tree process

Extra

The mean number of elements of a group at each step t is given by

$$f(x_t, t) = x_t - \sum_{j=0}^{x_t-1} \left(\sum_{i=0}^j \binom{x_t}{i} \Pi_t^i (1 - \Pi_t)^{x_t-i} \right)^{S-t}$$

where Π_t is obtained with the hypergeometric distribution and considering that if a symbol has been used as *the most common* then it cannot be reused, and so at each step $S \rightarrow S - 1$ and $k \rightarrow k - 1$.

$$\Pi_t = \frac{\binom{S-1-t}{k-1-t}}{\binom{S-t}{k-t}} = \frac{\Gamma(S-t)\Gamma(k-t+1)}{\Gamma(k-t)\Gamma(S-t+1)}$$

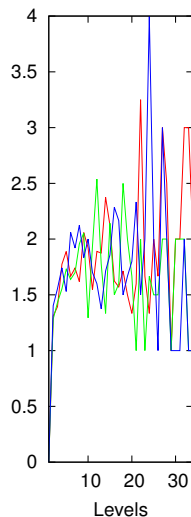
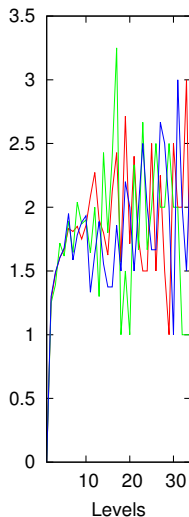
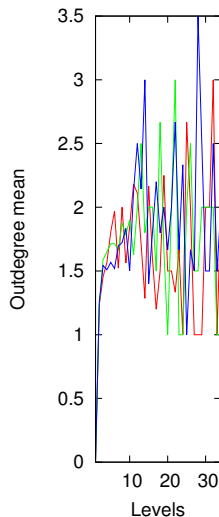
The process for the mean number of elements can so be defined as

$$x_{t+1} = f(x_t, t)$$

where $x_0 = x(0)$ and is equal to \mathcal{N} for the main process.

Null model 1

Extra



Null model 2

Extra

