# Politechnika Wrocławska

# Towards identifying software project clusters with regard to defect prediction

Marian Jureczko, Wrocław University of Technology
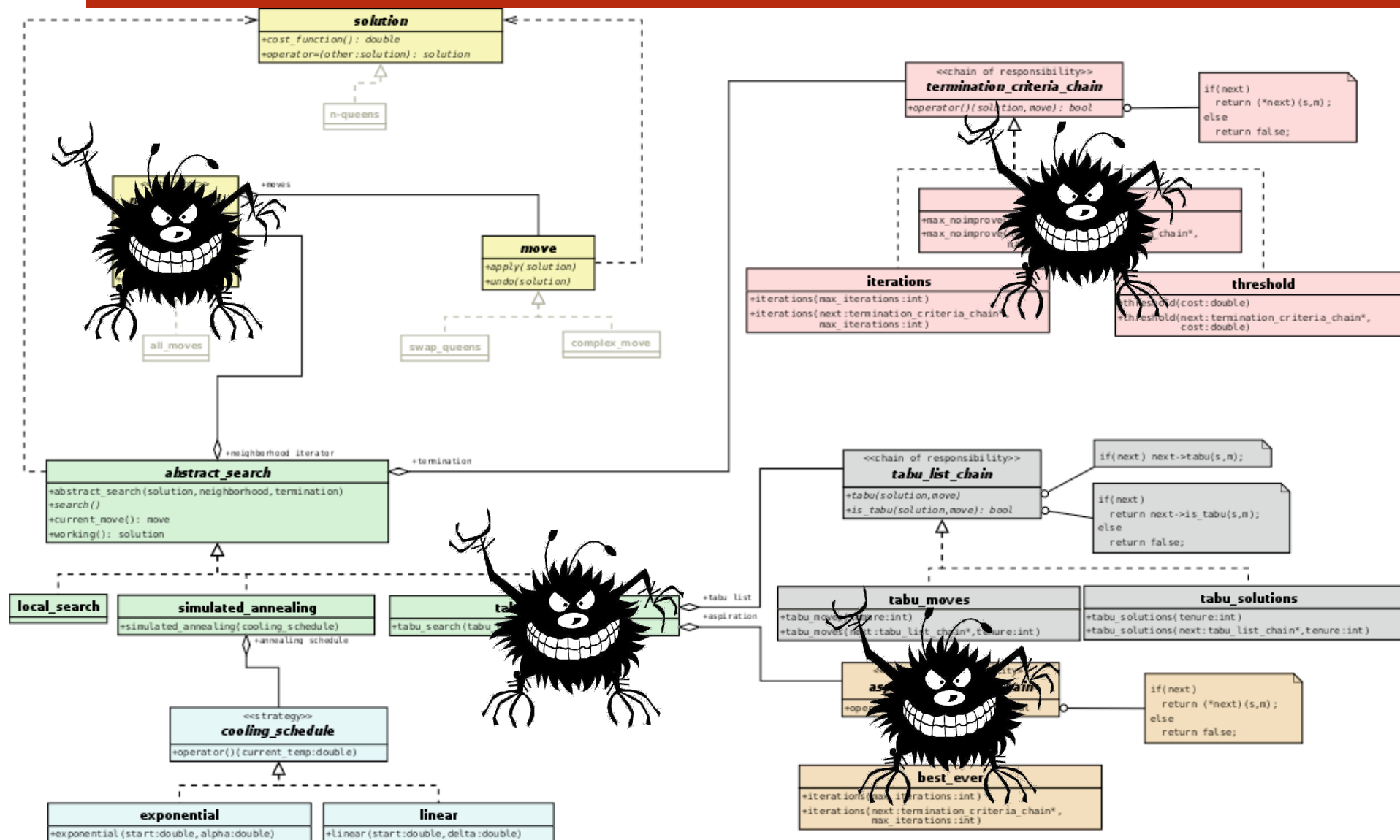Lech Madeyski, Wrocław University of Technology

# Agenda

- Introduction

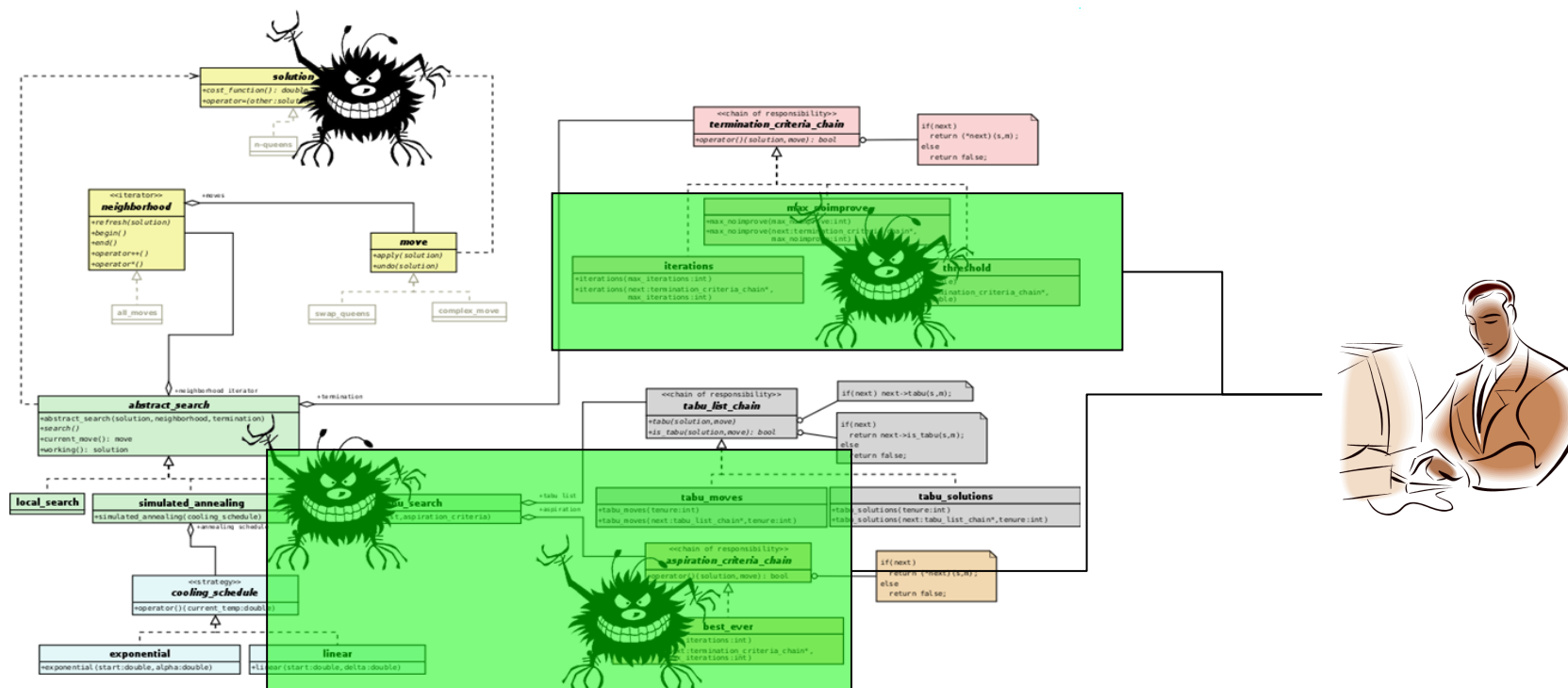- Data acquisition

- Study design

- Results

- Conclusions

# Introduction

**solution**

+cost_function(): double
+operator=(other:solution): solution

n-queens

+moves

**move**

+apply(solution)
+undo(solution)

all_moves

swap_queens

complex_move

<<chain of responsibility>>
**termination_criteria_chain**

+operator()(solution,move): bool

```
if(next)
    return (*next)(s,m);
else
    return false;
```

+max_noimprove
+max_noimprove

**iterations**

+iterations(max_iterations:int)
+iterations(next:termination_criteria_chain,
            max_iterations:int)

**threshold**

+threshold(cost:double)
+threshold(next:termination_criteria_chain,
           cost:double)

+neighborhood iterator

+termination

**abstract_search**

+abstract_search(solution,neighborhood,termination)
+search()
+current_move(): move
+working(): solution

<<chain of responsibility>>
**tabu_list_chain**

+tabu(solution,move)
+is_tabu(solution,move): bool

```
if(next) next->tabu(s,m);
```

```
if(next)
    return next->is_tabu(s,m);
else
    return false;
```

local_search

**simulated_annealing**

+simulated_annealing(cooling_schedule)

**tabu**

+tabu_search(tabu

+tabu list

+aspiration

**tabu_moves**

+tabu_moves(tenure:int)
+tabu_moves(next:tabu_list_chain*,tenure:int)

**tabu_solutions**

+tabu_solutions(tenure:int)
+tabu_solutions(next:tabu_list_chain*,tenure:int)

+annealing schedule

<<strategy>>
**cooling_schedule**

+operator()(current_temp:double)

**as...chain**

+oper...

```
if(next)
    return (*next)(s,m);
else
    return false;
```

**best_ever**

+iterations(max_iterations:int)
+iterations(next:termination_criteria_chain,
            max_iterations:int)

**exponential**

+exponential(start:double,alpha:double)

**linear**

+linear(start:double,delta:double)

# Motivation – Why defect prediction?

20% of classes contain 80% of defects



We can use the software metrics to predict error prone classes and therefore prioritize and optimize tests.

# Motivation – Why clustering projects?

- Defect prediction is sometime impossible because lack of training data:
  - It may be the first release of a project
  - The company or the project may be to small to afford collecting training data

- With well defined project clusters the cross-project defect prediction will be possible

# Definitions

- Defect
  - Interpreted as a defect in the investigated project
  - Commented in the version control system (CVS or SVN)
- Defect prediction model

Values of Metrics
for a given
java class

- WMC = ...
- DIT = ...
- NOC = ...
- CBO = ...
- RFC = …
- LCOM = …
- Ca=...
- ....

Model

Estimated
Number
of
Defects

# Data acquisition

- 19 different metrics were calculated with the CKJM tool (http://gromit.iiar.pwr.wroc.pl/p_inf/ckjm)
  - Chidamber & Kemerer metrics suite
  - QMOOD metrics suite
  - Tang, Kao and Chen's metrics (C&K quality oriented extension)
  - Cyclomatic Complexity, LCOM3, Ca, Ce and LOC
- Defects were collected with BugInfo ( http://kenai.com/projects/buginfo)

# Data acquisition

- 92 versions of 38 projects were analysed
  - 6 proprietary projects *(5 custom build solutions from insurance domain, 1 quality assurance tool)*
  - 17 academic projects
  - 15 open-source projects *(Apache Ant, Apache Camel, Ckjm, Apache Forrest, Apache Ivy, JEdit, Apache Log4j, Apache Lucene, PBeans, Apache POI, Apache Synapse, Apache Tomcat, Apache Velocity, Apache Xalan-Java, Apache Xerces)*
- Metrics Repository ( http://purl.org/MarianJureczko/MetricsRepo )

Politechnika Wrocławska

# Study design - clustering

# Results

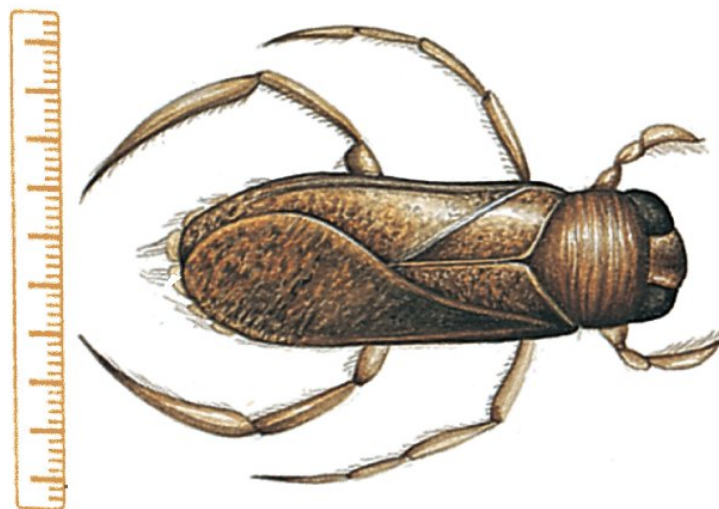| Cluster | Is the cluster model better? | P value (statistical test) |
|---|---|---|
| 1st of 2 | YES | 0.954 |
| 2nd of 2 | NO | - |
| proprietary A | NO | - |
| proprietary B | YES | 0.035 |
| proprietary / open | YES | 0.005 |
| open-source | NO | - |

# Results

- Cluster 'Proprietary B'
    - custom build solutions;
    - heavy weight, plan driven development process;
    - already installed in the customer environment;
    - insurance domain;
    - manual tests;
    - similar development period;
    - use database;
    - proprietary – the same company.

- Cluster 'proprietary / open'
    - text processing domain;
    - SVN and Jira or Bugzilla used;
    - medium size international team;
    - automatization in the testing process;
    - do not use database

# Conclusions

- 92 releases of 38 proprietary, open-source and academic projects were analysed

- 2 methods of clustering were applied

- 6 clusters were identified and the existence of 2 of them were proven

Politechnika Wrocławska

**Thank You
for Your attention**