**CHALMERS** UNIVERSITY OF TECHNOLOGY | UNIVERSITY OF GOTHENBURG

MRV Chaudron Hamburg 2019

# Empirical Studies into Modelling in Software Development in the Age of Big Data and AI

**Michel R.V. Chaudron**
Professor in Software Engineering
Joint Computer Science and Engineering dept.
Chalmers and Gotenburg university, Sweden
chaudron@chalmers.se

---

**CHALMERS** UNIVERSITY OF TECHNOLOGY | UNIVERSITY OF GOTHENBURG

MRV Chaudron Hamburg 2019

# Outline of talk

**Introduction**

- **Modeling in Software Development**
    - **Motivation**
- **Description of Practice of Modeling**

**Empirical Research in Software Design and Modeling**

- **Topics include:**
    - **illustrations of various of modeling-related 'big' data**
    - **illustrations of machine learning in this field**

**Summary & Conclusions**

Questions are welcome

**CHALMERS** | UNIVERSITY OF GOTHENBURG

# Joint Dept of Computer Science
# Chalmers and Gotenborg Univ`s: Divisions

Computing Science

Information Security

Formal Methods

Interaction Design

Functional Programming

Software Engineering

Networks and systems
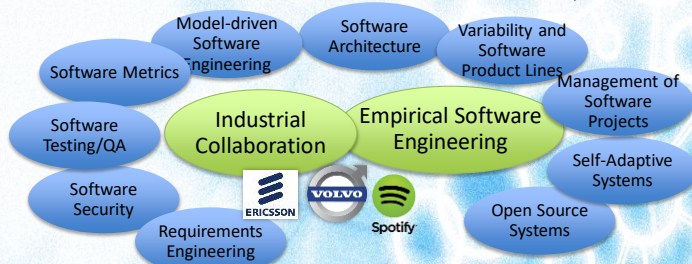
Computer Engineering

4



**CHALMERS** | UNIVERSITY OF GOTHENBURG

# Software Engineering Division Staff

Miroslaw Staron

Richard Berntsson

Robert Feldt

Michel Chaudron

Jan Bosch

Ivica Crnkovic

Gregory Gay

Richard Torkar

Jennifer Horkoff

Patrizio Pelliccione

Francisco Oliveira Neto

Software Metrics

Model-driven Software Engineering

Software Architecture

Variability and Software Product Lines

Management of Software Projects

Software Testing/QA

Industrial Collaboration

Empirical Software Engineering

Self-Adaptive Systems

Software Security

Requirements Engineering

Open Source Systems

Thorsten Berger

Christian Berger

ERICSSON   VOLVO   Spotify

Agneta Nilsson

Riccardo Scandariato

Gul Calikli

Jan-Philipp Steghöfer

Eric Knauss

Philip Leitner

Regina Hebig

Birgit Penzenstadler

---

**CHALMERS** | UNIVERSITY OF GOTHENBURG

# SOFTWARE ENGINEERING CONFERENCES IN GOTHENBURG

Open Source 2016
REFSQ 2016
ICSA 2017    http://icsa-conferencesorg/2017/
Mensura 2017    http://www.iwsm-mensura.org/2017
IFIPTM 2017 Int. Conf. on Trust Management
     http://wp.portal.chalmers.se/ifiptm2017/

**ICSE** 40TH INTERNATIONAL CONFERENCE ON SOFTWARE ENGINEERING

+ yearly Lindholmen Software Development day (600+)

---

**CHALMERS** | UNIVERSITY OF GOTHENBURG

## Software Engineering in Africa (SEiA)

- Software Engineering in Africa workshop at ICSE
- SIDA BRIGHT project (2015-2010)
  - Supervise 10 lecturers in Uganda towards their PhD degree.
  - Yearly summerschools in East-Africa



- Tentative plans:
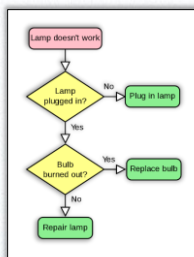  African Summerschool on Software Engineering 2020

Modelling History

Juha-Pekka Tolvanen, CEO Meta-Case, Keynote at Code Generation 2014:
The business cases of modeling and generators



Example 4+1 Views model

## Fake UML News – Urban Myths

UML is a notation, not a method –

UML can only be used with object orientation

UML is not a process
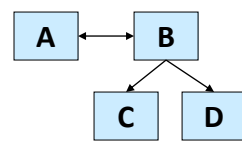> i.e. UML is not tied to waterfall or RUP or SCRUM or …

Only when used for code generation, is modeling effective

Michel Chaudron is a UML advocate

BUSTED

---

# The Importance of Architecture

Cost of SW fault repair costs

Architecting/design determines 90% of the costs and risks

≥ 100x

Analyse before you build.

1x

Requirements & Architecting

Design

Implementation

Testing

Deployment

Operation

Increasing Complexity of Software

Context factors
- Complex system with many modules and dependencies

'Management Design Complexity'

Modeling is inevitable



Motivation

Why should we care about modeling software design?

Pro-Modeling-camp:
"I can not imagine we could develop our systems without any modeling."

Anti-Modeling-camp:
"Models are useless. Code is the only truth."

01/04/2019

---

# Modeling & Design in Civil Engineering



Context factors
- Complex endeavor
- Large project team
- Multidisciplinary expertise
- Large scope over time

*'Knowledge Management', Planning & Coordination*

---

# Summary of Arguments in favour of Modelling in Software Development

+ Communicating / Coordination

+ Analysing / Predicting

+ Understanding / Structuring

+ Guiding Construction

+ Blueprint for Production

01/04/2019





9

CHALMERS | UNIVERSITY OF GOTHENBURG

# Not always rational …



---

CHALMERS | UNIVERSITY OF GOTHENBURG

# Modeling & Documentation
# in Agile Development

- Agile principles:
  *working software over comprehensive documentation*



How do you feel about documentation at work?

Survey under
75+ agile developers

Tending towards:
 we need a bit more

Modeling is
compatible with
agile development

Christoph J. Stettina and Werner Heijstek, Necessary and Neglected? An Empirical Study of Internal Documentation in Agile Software Development Teams 29th ACM Int. Conf. on Design of Communication , Pisa, Italy

CHALMERS | UNIVERSITY OF GOTHENBURG

MRV Chaudron Hamburg 2019

# Informal (whiteboard) drawing

Easy to make.

But not so easy to:
- Change / update
- Search
- Link / Trace
- Analyse / Check

CHALMERS | UNIVERSITY OF GOTHENBURG

# To Model or Not to Model?

CHALMERS | UNIVERSITY OF GOTHENBURG

MRV Chaudron Hamburg 2019

# Better Questions

- Which people in my development organization benefit from design models? Through which task/activity?
- Which information should I include in my models?
  - Who 'owns' this information?
- How can we integrate modeling into our
  - development process?
    - who should be responsible for creating and maintaining models?
  - development toolchain?

---

CHALMERS | UNIVERSITY OF GOTHENBURG

MRV Chaudron Hamburg 2019

# Understanding the practice of modeling in software projects

CHALMERS | UNIVERSITY OF GOTHENBURG

# The Model-based SE Spectrum

| Code only | Design for guiding | Design for guiding & update | Model-centric | Model only |
|---|---|---|---|---|
| | Model A/D/I | Model D/I | Model I | Model ? |
| Code | Code | Code | Code | |

**Fig. 1.** The modeling spectrum

Model driven architecture: Principles and practice, Brown, A.W., Software and Systems Modeling, 2004

---

CHALMERS | UNIVERSITY OF GOTHENBURG

# Styles of UML modeling

**Sketch**

**Communication**

More effort ⇒ More expensive

**Recipe for construction**

|  | Ideation | Externalization | Production |
|---|---|---|---|
| Building | | | |
| Animation | | | |
| Software | | | |

# Evolution of Models during a project

|  | Ideation | Externalization | Production |
|---|---|---|---|
| Architecture | | | |
| Design | | | |
| Implementation | | | |

Different stages have different needs on formality
Yet tasks are connected:
        'smooth' handover over design representation is needed

# Target audience of the models?

| Ideation | Externalization | Production |
|----------|-----------------|------------|

Main audience:        **people**                    **computers**

Ok with informal spec
(has context, domain
& background knowledge)

Needs
formal and
complete spec

---

# Empirical Research in Modeling in Software Design

The challenges with Industrial ('Real') Data

1. Difficult to get access to
2. Difficult to publish / share (replicate)

DIY - Let's Build Our Own Dataset

Können wir das schaffen?!



2 main pillars for Big Data

Collecting data

Analysing data

That looks doable

CHALMERS
UNIVERSITY OF TECHNOLOGY | UNIVERSITY OF GOTHENBURG

# Collecting Big Data

Bilal
Karasneh

Idea: This is a 'Search'-problem.

Probably Google can do this for us.

Attempt 0.1 & 0.2:

1. Search for (filetype) .xmi using Google

   Unsuccessful: incredibly much that is not about 'XMI' for UML

   Manual filtering is very time-consuming

2. Search for 'class diagram' via Google Image Search

   + : many images can be found, not all are class diagrams

   - : very high variation in quality
       tutorial/lecture-slides (on notation), student-projects, …

---

CHALMERS
UNIVERSITY OF TECHNOLOGY | UNIVERSITY OF GOTHENBURG

# Google Image Results – page 1

CHALMERS | UNIVERSITY OF GOTHENBURG

# Google Image Results – **page 2**



Many many false positives.

CHALMERS | UNIVERSITY OF GOTHENBURG

# Google Image Results – **page 3**



Still many false positives.

After much manual filtering: collection of about 1000 diagrams.

---

# Classifier for recognizing UML Class diagrams

With my students:
- Ingimar Samúelsson
- Jóel Hjaltason

- Input:
  arbitrary image (.jpeg)
- Output:
  Yes/No (=UML Class Diagram)
- Initial filtering done by rules
  – Size (no icons), colours

- 1300 Diagrams, of which
  650 UML CD and 650 other
  (incl. Seq.D, ER, charts, …)

| Features include: |
| --- |
| Contours & Shapes |
| Horizontal and Vertical Lines |
| Rectangle-H-W-ratio |
| Connected rectangles |
| Rectangles with 'divider-line' |
| NON-Hor/Ver lines |
| % of area used |

| Algorithm | Precision |
| --- | --- |
| Random Forrest | 0.95 |

---

# Automated Extraction of Models from Images

Bilal Karasneh



Recognizing:
- Rectangles
- Lines
- Text,
- arrow-heads

Input: .jpeg, .png

Output: XMI for the class model

**Milestone: IMG2UML Tool**

01/04/2019

---

# Is Reverse Engineering a solution?

- Automatically generated from source code, hence always up-to-date!
- But, …

---

CHALMERS | UNIVERSITY OF GOTHENBURG
UNIVERSITY OF TECHNOLOGY

# Reverse Engineering of a small system



Clearly different from forward designed UML designs
(o.a. in size, layout, detail, naming, ….)

---

# Classifier for UML Class Diagram Styles using Machine Learning

- Forward vs Reverse Engineered Diagram Classifier
- 16 Features
  - mostly related to 'parameters' of methods in diagram
- Sample of 999 class diagrams
  - Reverse : 806
  - Forward : 193

| Algorithm | Accuracy | Precision | Recall | F-measure |
|-----------|----------|-----------|--------|-----------|
| Random Forrest | 90.74 | 0.95 | 0.93 | 0.94 |

---

Collecting data

# The Quest for UML in Open Source Projects
## -- Mining GitHub --

**Regina Hebig** [*]    **Truong Ho-Quang** [*]    **Miguel-Angel Fernandez** [+]    **Gregorio Robles** [+]    **Michel R.V. Chaudron** [*]

[*] Gothenburg and Chalmers University
[+] Universidad Rey Juan Carlos

Hebig, Regina, et al. "The quest for open source projects that use UML: mining GitHub." *Proceedings of MODELS* 2016.

---

**Slide 1**

CHALMERS | UNIVERSITY OF GOTHENBURG

Collecting data

# Some statistics

July 2015

GHTorrent    GitHub

12 850 000 non-forked repos

① Data collection

Potential UML file list

② Filter UML files
- UML Image Filter
- Textual Filter
- Validation

UML File list

③ Extract Meta-data

CVSAnalY MySQL

④ Query

⑤ Analyse result

December 2016

Identified
93 000+ UML files
24 000+ repositories

---

**Slide 2**

CHALMERS | UNIVERSITY OF GOTHENBURG

Collecting data

## Surprise 1
## UML models appear in many formats

| File type | .uml | xmi | svg | png | jpeg | bmp | jpg | gif |
|-----------|------|-----|-----|-----|------|-----|-----|-----|
| **Share** | 34% | 4% | 1% | 37% | <1% | <1% | 10% | 13% |

image formats

60%

NB. We did not *yet* look into: pdf, word, ppt, specific UML-CASE tooling-formats (.uml, .umple, …)

---

## Projects grouped by number of UML files

| #UML files | 1 | [2-9] | [10-99] | [100,∞) |
|---|---|---|---|---|
| #Repositories | 14 749 | 8 567 | 1 333 | 68 |

Projects grouped by #UML files



## Example analysis: size of class diagrams

Distribution of diagrams by number of classes

CHALMERS | UNIVERSITY OF GOTHENBURG

# Countries of Projects that use UML

Distribution of participant by continents

- Europe
- North America
- South America
- Asia
- unknown
- Oceania
- Africa

54%
18%
11%
12%
3% 1% 1%



CHALMERS | UNIVERSITY OF GOTHENBURG

# Opportunity

## This dataset enables many types of empirical studies

### Are you looking for models?

We have collected    http://oss.models-db.com/

**35 000+** class diagrams

from **24 000+** open source projects at GitHub

These diagrams can be traced back to the projects, hence it is possible to find associated project data such as source code, commit messages, commit-dates, and much more.

We would love to hear from you:
- Would you like to use the dataset in your research?
- Which research questions do you recommend us to look at?

Hebig, R. & Ho-Quang, T. & Robles, G. & Fernandez, M.A. & Chaudron, M.R.V. (2016). The Quest for Open Source Projects that Use UML: Mining GitHub In proceedings, *ACM/IEEE 19th International Conference on Model Driven Engineering Languages and Systems*, Saint-Malo, France, October 2-7, 2016.

Dataset

http://oss.models-db.com/

Research Group

Regina Hebig    Truong Ho-Quang    Miguel-Angel Fernandez

Gregorio Robles    Michel R.V. Chaudron

Contact person:    Truong Ho-Quang    truongh@chalmers.se

CHALMERS | UNIVERSITY OF GOTHENBURG    Universidad Rey Juan Carlos

## Research Challenge

tools

Value/Utility

use

discovery

Knowledge

Patterns?
Which ones to look for?

Data

+/- 100,000

Enrich via labelling:
- Reverse vs Fwd design?
- Student project?
- Domain model vs
  design model?

---

## Knowledge Discovery Challenge

- Huge variety of data:
  - Graphs: UML diagrams
  - Source code
  - Text: SAD, bug-reports, ….

Currently looking beyond models and code, also into Software Architecture/Design Documents

Approach to Modeling is contextual



Is there a DNA for Software Design?

"Sequence motifs are short recurring patterns in DNA that are presumed to have biological function."

D'haeseleer, P. *Nature Biotechnology* **24**, 423 - 425 (2006).

Image taken from bio.miami.edu

## Slide 1

CHALMERS UNIVERSITY OF TECHNOLOGY | UNIVERSITY OF GOTHENBURG

# Challenge: Uncovering the hidden structure of designs

Rebecca Wirfs-Brock: Software designs are built from building blocks that have stereotypical responsibility-roles

| Controller | Coordinator | Interfacer |
| Structurer | Service Provider | Information Holder |

## Slide 2

CHALMERS UNIVERSITY OF TECHNOLOGY | UNIVERSITY OF GOTHENBURG

# Definition of Role Stereotypes

There are many ways to understand the nature of a class, but I start by looking at its name. *Aptronyms* are names that match a person's occupation—Joe Strong is a weight lifter; Suzie Snow is a ski instructor. Because I look at a class's name to suggest its role in a design, I expect class names to be aptronyms. For example, in Java, a StringTokenizer picks apart segments of a string, and the ClassLoader loads classes. But names aren't always illuminating. So I also scan a class for intention-revealing method names that suggest the class's behavior. Of course, the definitive source is always the code, but I shouldn't have to pore over details just to get the gist of a class.

In this column, I introduce several characteristics I ascribe to classes when trying to understand their nature and purpose. I hope you find these useful quiddities and not mere quibbles. (I'm both delighted by and distrustful of a word that has definitions with opposite meanings. The technical term for such a beast is autoantonym.)
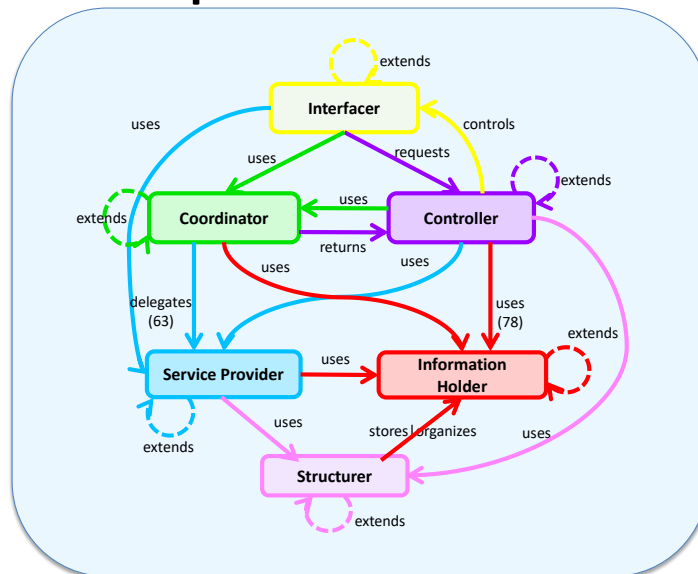
**Role stereotypes**

Purposeful oversimplifications, or *role stereotypes* from Responsibility-Driven Design

- *Information holder*: an object designed to know certain information and provide that information to other objects.
- *Structurer*: an object that maintains relationships between objects and information about those relationships. Complex structurers might pool, collect, and maintain groups of many objects; simpler structurers maintain relationships between a few objects. An example of a generic structurer is a Java HashMap, which relates keys to values.
- *Service provider*: an object that performs specific work and offers services to others on demand.
- *Controller*: an object designed to make decisions and control a complex task.
- *Coordinator*: an object that doesn't make many decisions but, in a rote or mechanical way, delegates work to other objects. The Mediator pattern is one example.
- *Interfacer*: an object that transforms information or requests between distinct parts of a system. The edges of an application contain user-interfacer objects that interact with the user and external interfacer objects, which communicate with external systems. Interfacers also exist between subsystems. The Facade pattern is an example of a class designed to simplify interactions and limit clients' visibility of objects within a subsystem.

## Slide 1

CHALMERS | UNIVERSITY OF GOTHENBURG

# Relationships between role stereotypes



## Slide 2

CHALMERS | UNIVERSITY OF GOTHENBURG

### Automated Classification of Role-Stereotypes by Machine Learning

Arif Nurwidyantoro
Monash University
Melbourne, Australia
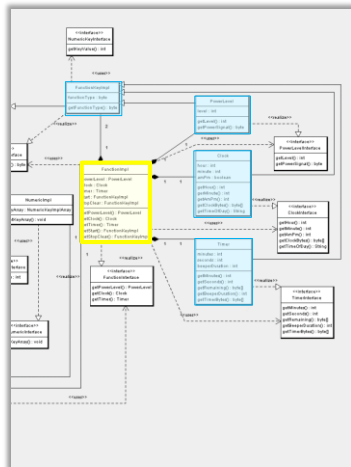arif.nurwidyantoro@monash.edu

Truong Ho-Quang
Chalmers | Gothenburg University
Gothenburg, Sweden
truongh@chalmers.se

Michel R.V. Chaudron
Chalmers | Gothenburg University
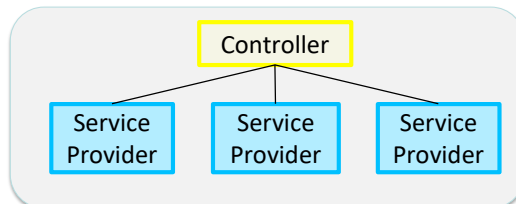Gothenburg, Sweden
chaudron@chalmers.se

- Build a ground truth for 3 open source projects
  - K9, BitCoinWallet, Home3D
- Features used:
  - source code (e.g. names), design metrics, network-metrics

# Mining graph-patterns
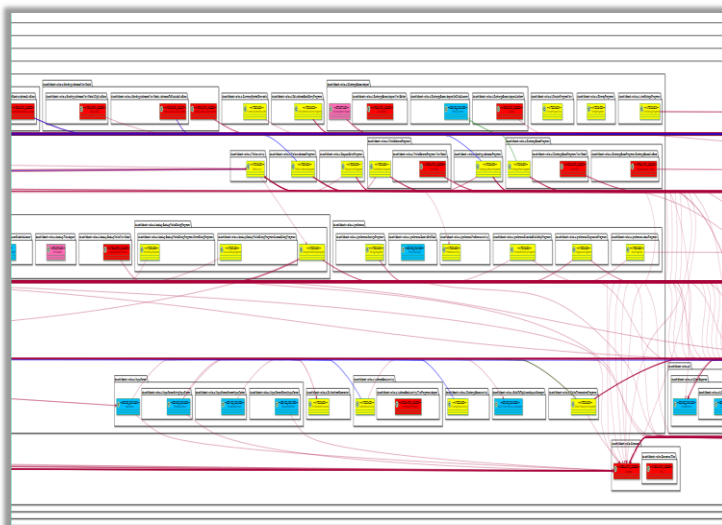# in Software Designs



Through labelling of roles, we find recurring patterns in the design



These patterns represent typical collaborations between responsibility-stereotypes
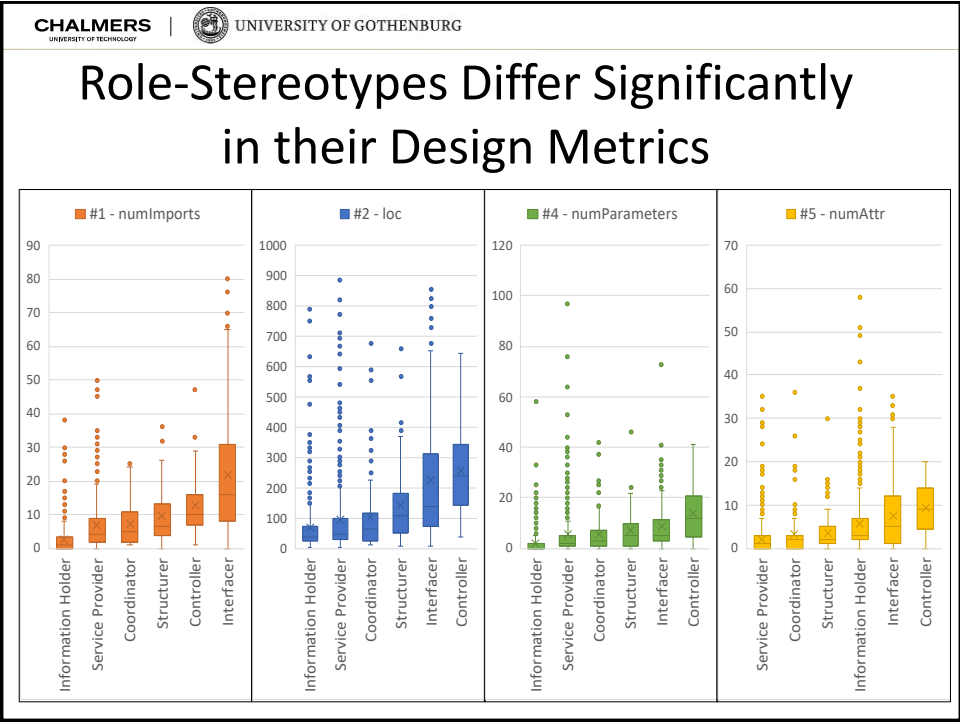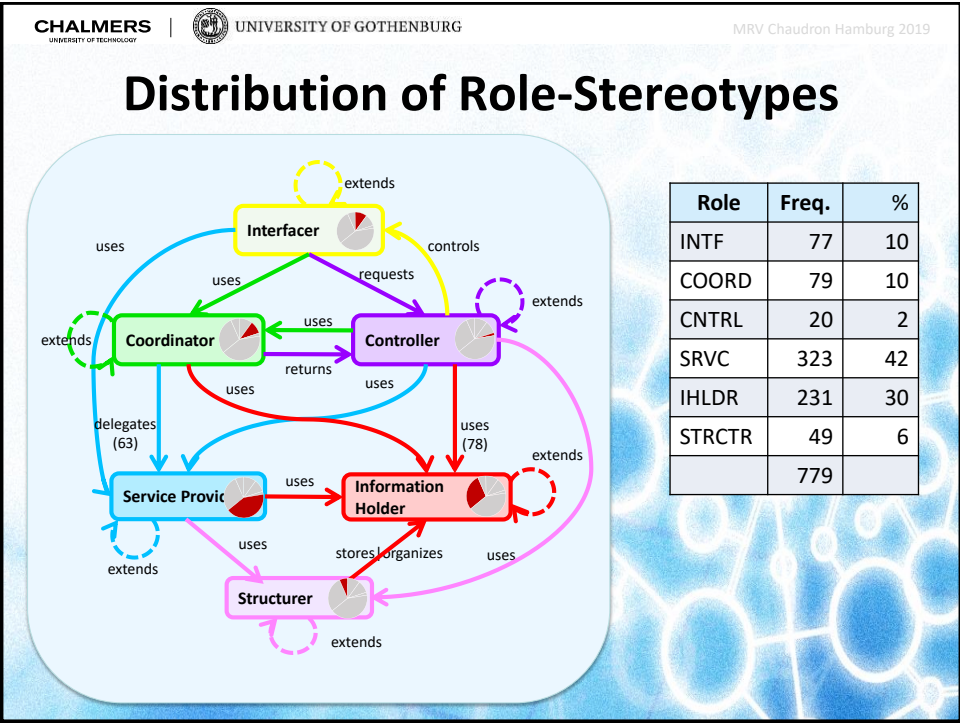
# Role-stereotypes in software design



Interfacer

Controller

Coordinator

Service Provider

Information Holder

Structurer

Distribution of Role-Stereotypes

| Role | Freq. | % |
|------|-------|---|
| INTF | 77 | 10 |
| COORD | 79 | 10 |
| CNTRL | 20 | 2 |
| SRVC | 323 | 42 |
| IHLDR | 231 | 30 |
| STRCTR | 49 | 6 |
| | 779 | |



Role-Stereotypes Differ Significantly in their Design Metrics

## Slide 1

**CHALMERS** | UNIVERSITY OF GOTHENBURG

# Role-stereotypes-based visualization of software design



- Interfacer
- Controller
- Coordinator
- Service Provider
- Information Holder
- Structurer

Ratio of different role-stereotypes within a component gives information about its main type of responsibility

## Slide 2

**CHALMERS** | UNIVERSITY OF GOTHENBURG    MRV Chaudron Hamburg 2019

# Possible uses of Role-Stereotypes

Role-Stereotypes aid in:
- Program understanding
- Design summarization
- QA for Architecture/Design
- Tailoring test-generation/coverage/…

## Challenge: Mining graph-patterns in Software Designs

Challenges:
- Stereotypes are 'idealized'.
  In practice, classes are not 'ideal' designs.

MRV Chaudron Hamburg 2019

## Knowledge Management in Software Development & Maintenance
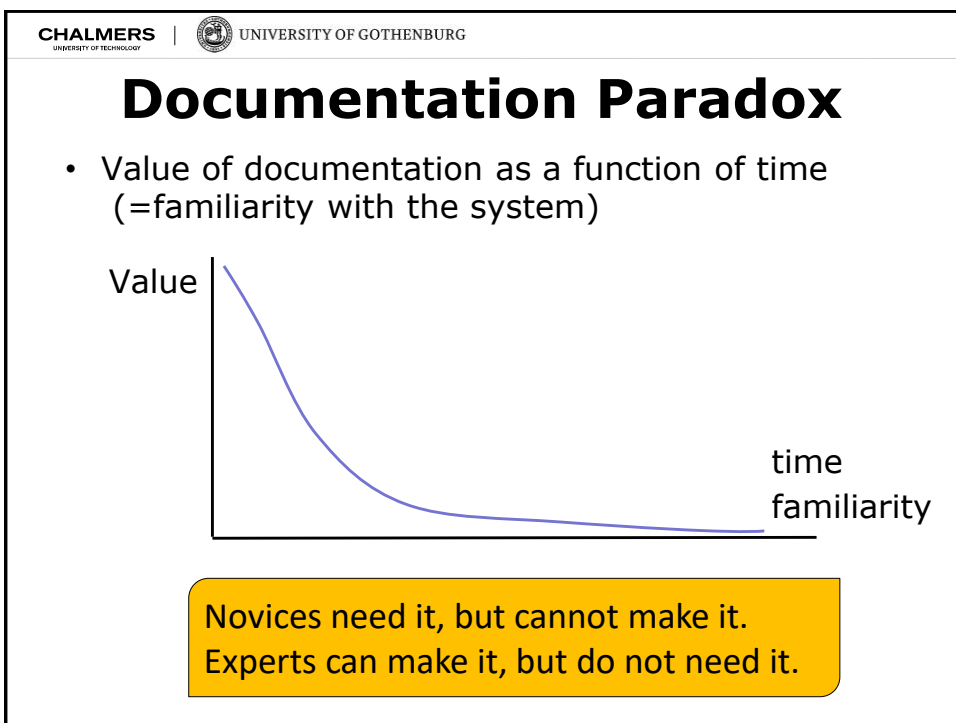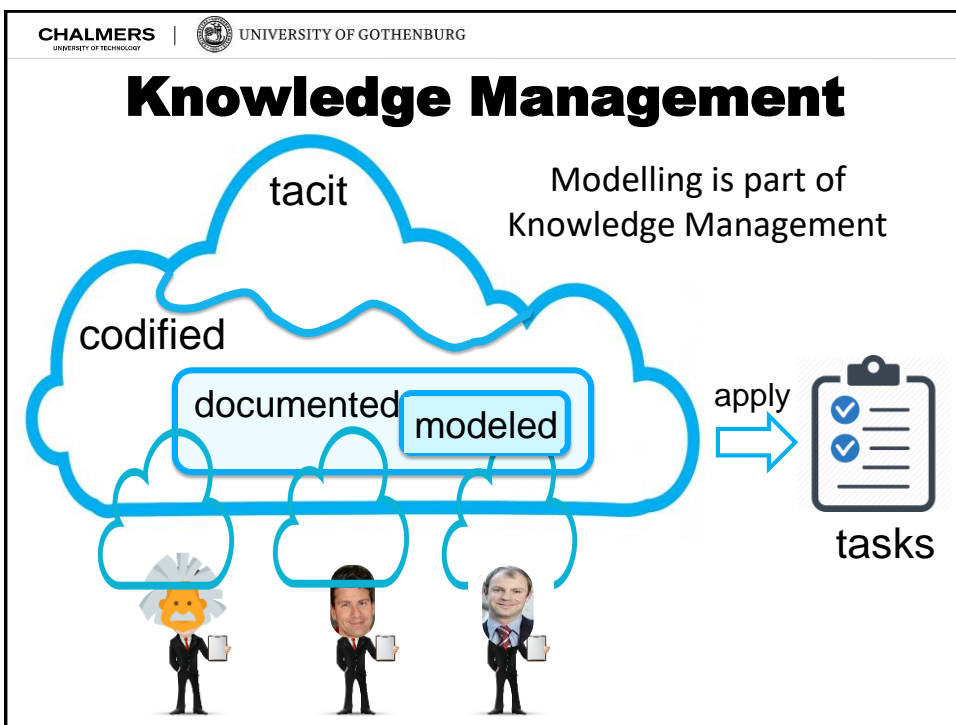
**Documentation Challenge:**
**Find & Navigate**

- Where is the information that I need?
- Many 'places'
  *GitHub, Project-PC, …*
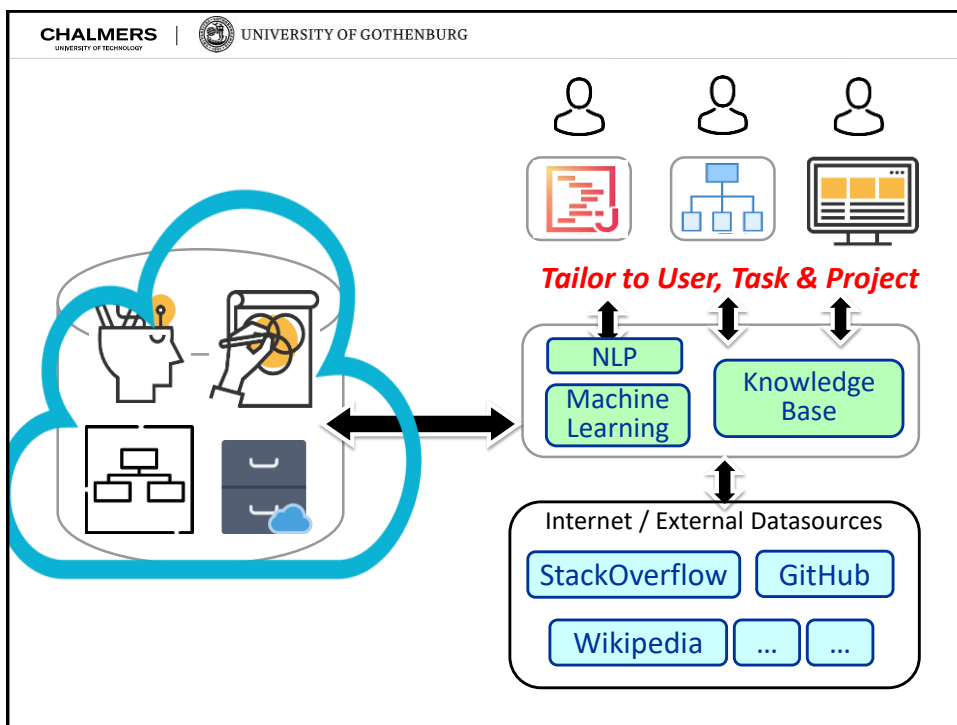- Poor searchability
- Poor navigability
- Is it up-to-date?



**Knowledge Management in**
**Software Development & Maintenance**

Reduces load on creating & maintaining documentation

- Extraction
  - NLP, image processing
- Linking
  - ontologies
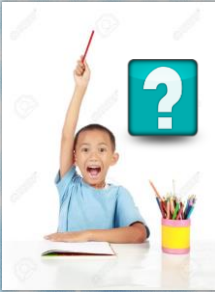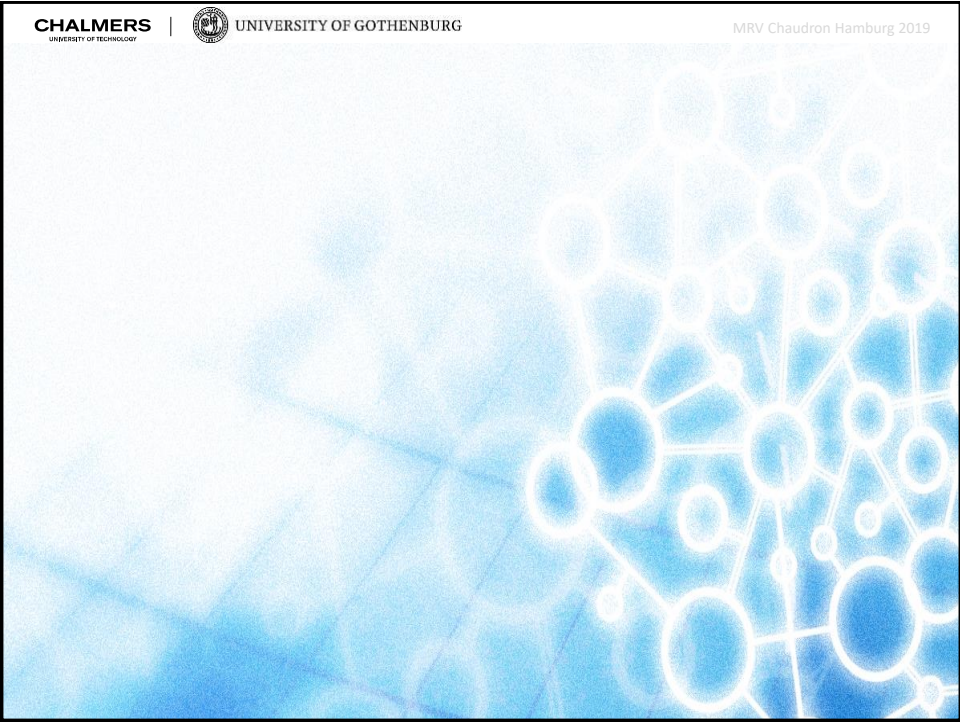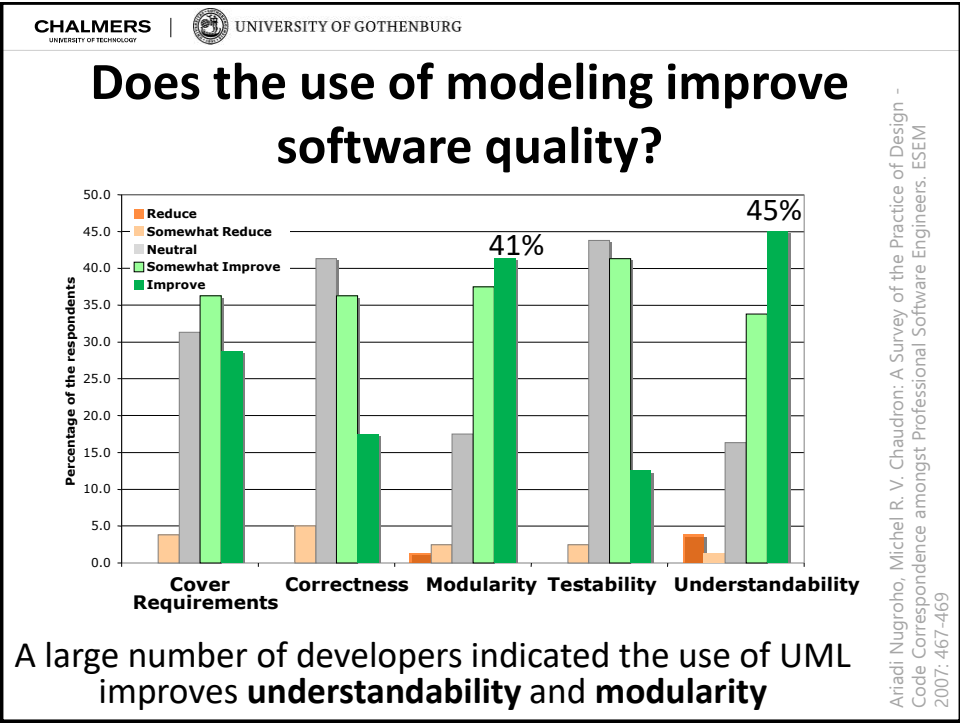- Up-dating
- Navigate / Present

## Concluding Remarks

- We have poor understanding of the consumption/use of models in practice
  - Mismatch between tools and tasks

- Engineering reality:
  - Modelling styles are manifold, driven by purpose & context
  - Lack of adoption of hygiene practices

- Next step / challenges
  - Knowledge-discovery
    - Heterogeneity of data
    - Absence of context / noise
  - Requires community effort

01/04/2019

**Does the use of modeling improve software quality?**

A large number of developers indicated the use of UML improves **understandability** and **modularity**

Ariadi Nugroho, Michel R. V. Chaudron: A Survey of the Practice of Design - Code Correspondence amongst Professional Software Engineers. ESEM 2007: 467-469

38

## Slide 1

# Summary of Evidence on Modeling

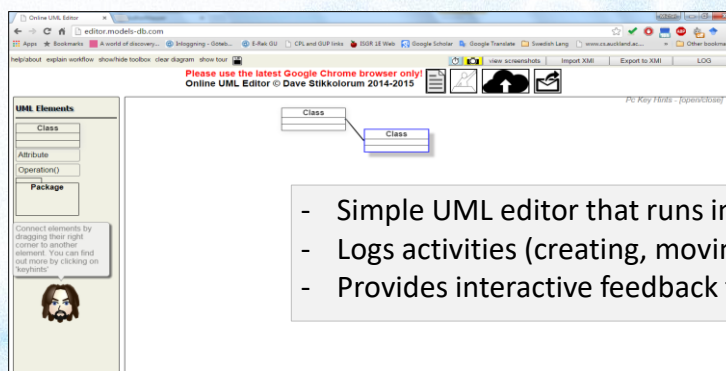| Benefit | Status |
| --- | --- |
| Communication | Confirmed |
| Analysis | No evidence |
| Structuring | No evidence |
| Guiding | Confirmed |

Main hurdles:
- migrating existing documentation
- Integrating modeling tools in toolchain
- Integrating modeling into process
- Keeping models & documentation up to date

## Slide 2

# Project : Towards an online learning environment for software design

with Dave Stikkolorum

- Simple UML editor that runs in a browser
- Logs activities (creating, moving, renaming)
- Provides interactive feedback to students

Goal is to integrate 'doing design' in online learning courses (UML/Analysis & Design/Softw architecture)

Fig. 2: LogViz visualisation tool

---

# Modelling style and model purpose



Styles of using UML
– as a sketch – thinking tool/understanding
– for communicating system design
– as a blueprint – guide the implementation work
– as a implementation (MDA) - code generation

C.F.J.Lange, M.R. V. Chaudron, J.Muskens. In Practice: UML Software Architecture and Design Description. IEEE Software 23(2): 40-46 (2006)

---

# Effect of Defects on Understandability of UML

ICSE 2005 paper
with Christian Lange

**Experiment:**
Show participants 2 types of UML models
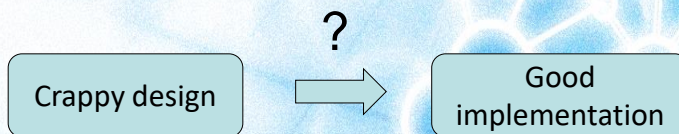A.  'good model'
B.  model omission / inconsistency

For 'faulty' models, participants more often differ in their interpretation.

Both students & professionals



Lange, Christian FJ, and Michel RV Chaudron. "Effects of defects in UML models: an experimental investigation." In *Proc 28th ICSE*, pp. 401-411. ACM, 2006.

# Challenge 1

Can we automatically assess the quality of a software design?
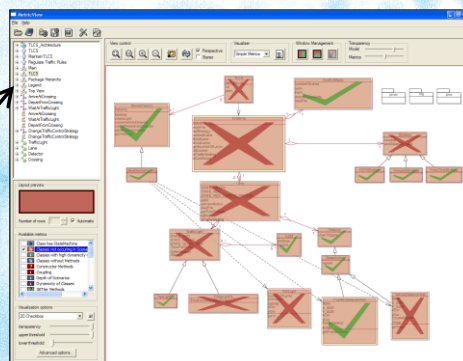
?

Crappy design ⟶ Good implementation

# Project: Quality Assessment of Software Design - the MetricView tool

When is a design 'good'?
Absence of 'bad things'?

with Christian Lange & Johan Muskens

Quality Metrics capture design principles
• Coupling
• Cohesion
• ... <extensible>



✗ = violation  ✓ = ok
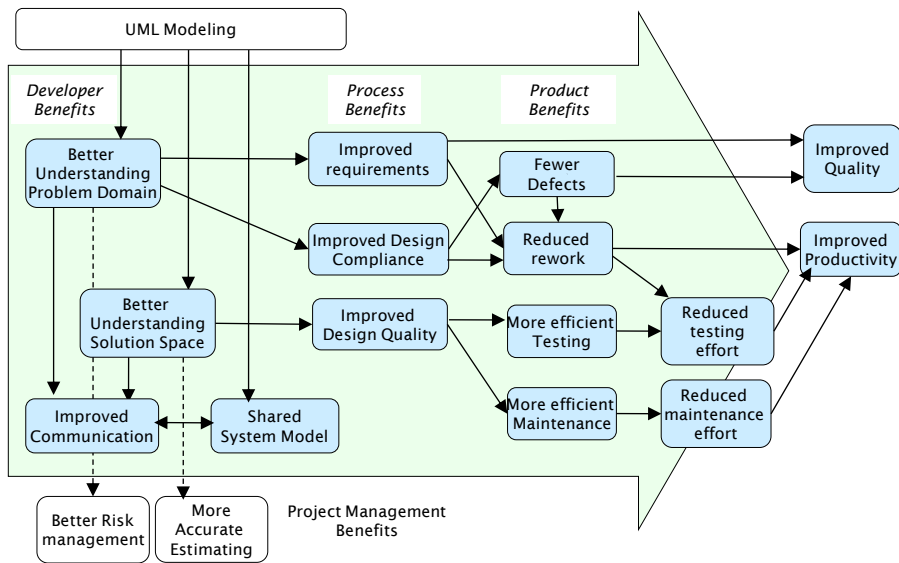
---

CHALMERS | UNIVERSITY OF GOTHENBURG

## Industrial Experience with MetricView as Quality Assurance Tool

- Based on 15+ industrial projects
- If there are weak spots in the design, then these are indicated as 'suspects' by MetricView
- About 90% of the weak spots indicated by MetricView do not require improvement according to project architects

> Syntactic checks (like metrics) are not sufficient to identify important areas in the design!

- The later MetricView is applied, the fewer 'weak spots' are removed from the design → process issue

## Theory of Benefits of Modeling



---

# Integrate Modeling into Process and Tooling

Naming and layout-conventions

Version Management

    Many tools around (e.g. CVS, SVN, …)

Reviews & Inspections

    Guidelines by e.g. Shull et.al., Biffl,

    Easy QA-tool: SDMetrics  http://www.sdmetrics.com/

Process

    Integrate into process:
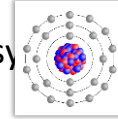
    Who? When? How? Update a model in documentation

    'Definition of done'

Low Hanging Fruit
Erin W 2010

# Definitions

**Model (noun)** =  an abstract representation
of a thing/sy[...]

often *systematic* representation

Models abstract: they focus on the
essential features and leave out others.

**Modeling** =  the process of creating a model
**(verb)**                 i.e. choosing what to represent
how to represent it

---

# Design (verb, noun)

Definition
**Design (v)** =   the process of making decisions
about something
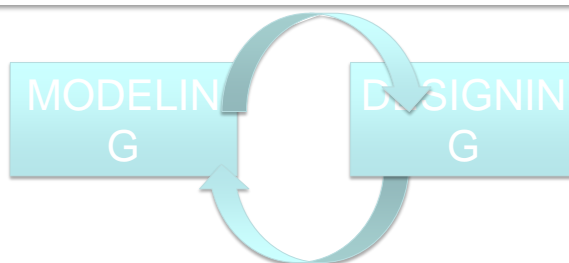that is to be                                                              built
or created**:**

**Design (n)**      =        the plans, drawings, etc., that
show how
something can be
made

Pitfall : 'design' & 'model' can be a verb and a noun
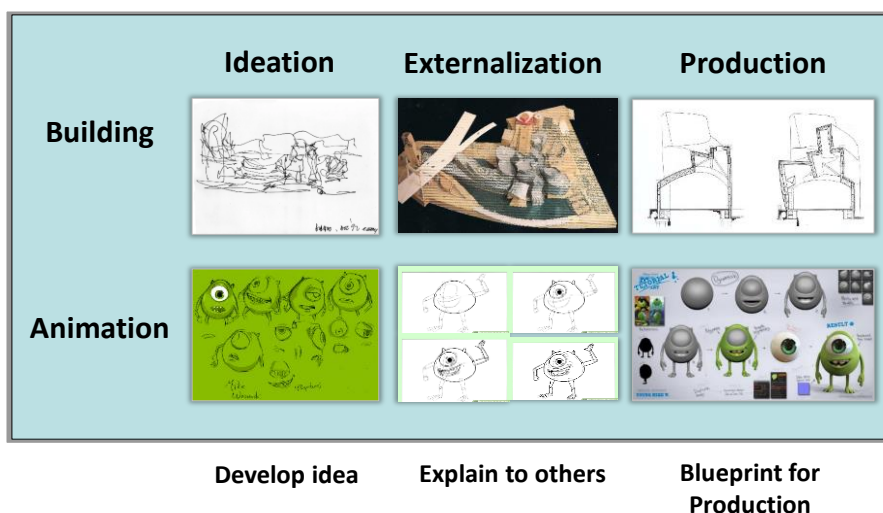
**CHALMERS** | UNIVERSITY OF GOTHENBURG

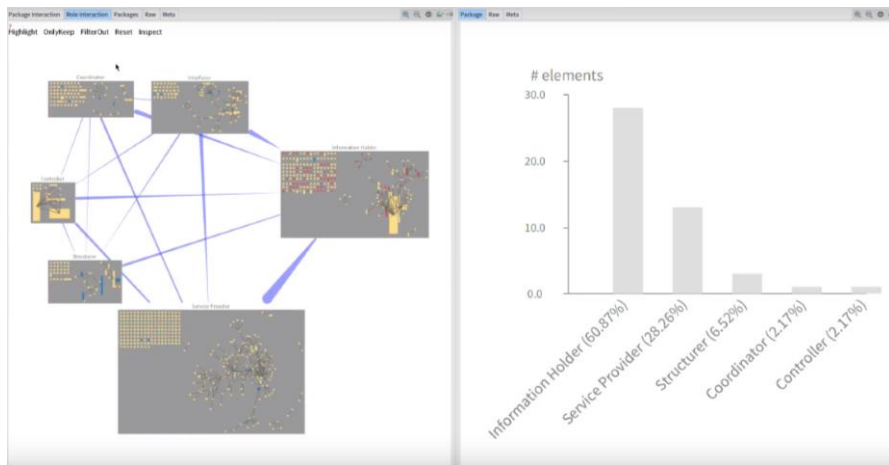# How do 'modeling' and 'design' relate?

Modeling and designing often go hand-in-hand:
   A model is used to understand, reason,    analyse, break-down, which leads to   adaptation and refinement of the design.

MODELING    DESIGNING

---

**CHALMERS** | UNIVERSITY OF GOTHENBURG

# Stages of Design & Modeling

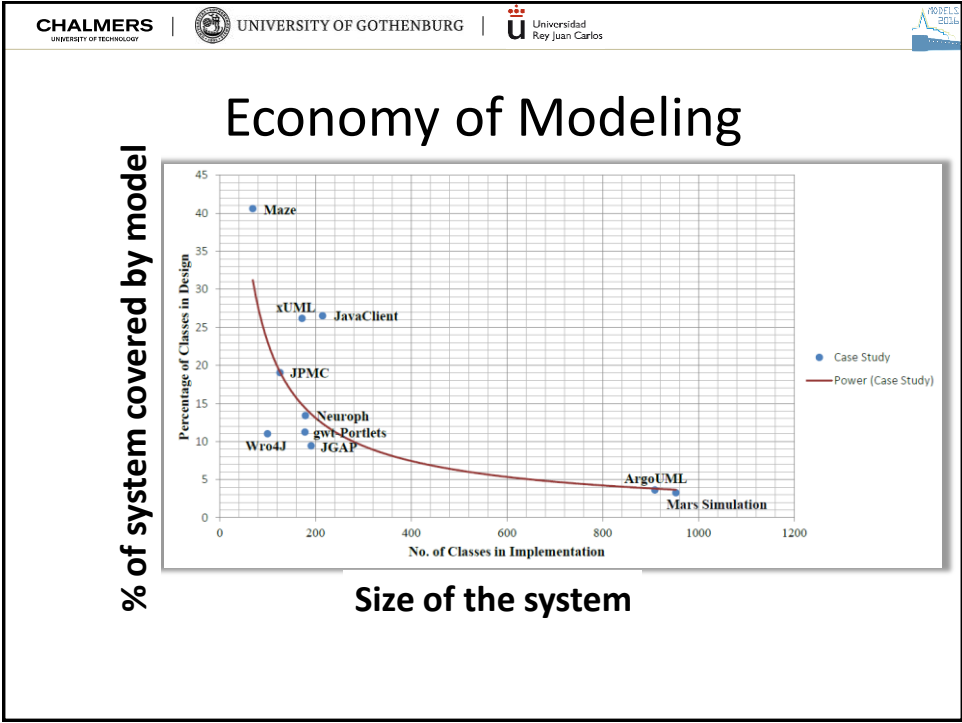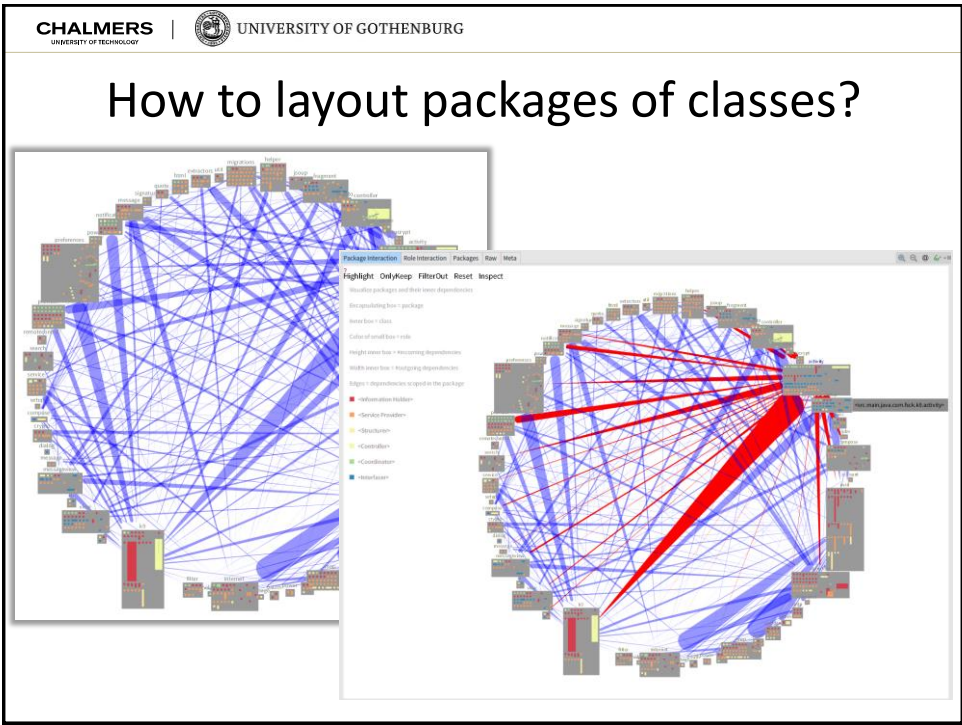|  | Ideation | Externalization | Production |
|---|---|---|---|
| **Building** | | | |
| **Animation** | | | |
|  | Develop idea | Explain to others | Blueprint for Production |

# Challenge: Visualisation



# Visualisation Challenge

- Scale
- Multiple uses
- Multiple abstraction levels
  - Models & Source Code

How to layout packages of classes?
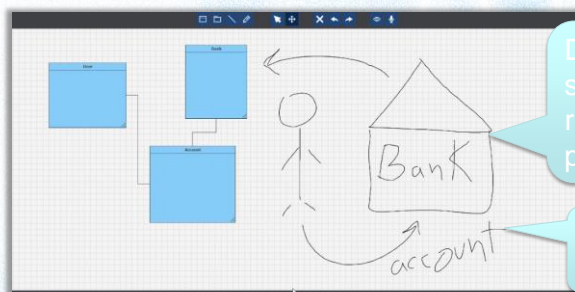


Economy of Modeling

**Integration of Tools in Design Flow and Development Process**

Are the tools the problem?

---



# Project : OctoUML
with Rodi Jolak

- Merge **whiteboard** (informal, free-form) and **CASE-tool** (formal) modeling

Digital editing for sketching: undo, resize, move, pan, zoom

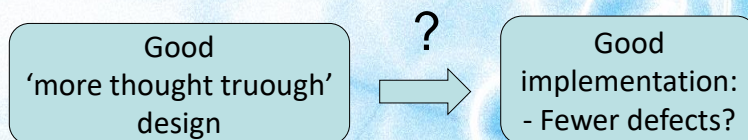Novel UI interaction models: touch, voice

https://www.youtube.com/watch?v=fsN3rfEAYHw&t=74s

Goal: Tools should support the *development process* (not only the modeling task)

# Challenge

Does the quality of a software design related to the quality of the implemented software?

| Good 'more thought truough' design | ? | → | Good implementation: - Fewer defects? |

---

# Project: Does Quality of Modeling Matter? An Industrial Case Study

with Ariadi Nugroho

Focus on *detail* in a UML Model

attributes

operations

associations

Low detail                    High detail

Relation between UML-LoD and Code Quality



Relation between Level of Detail and Defect Density