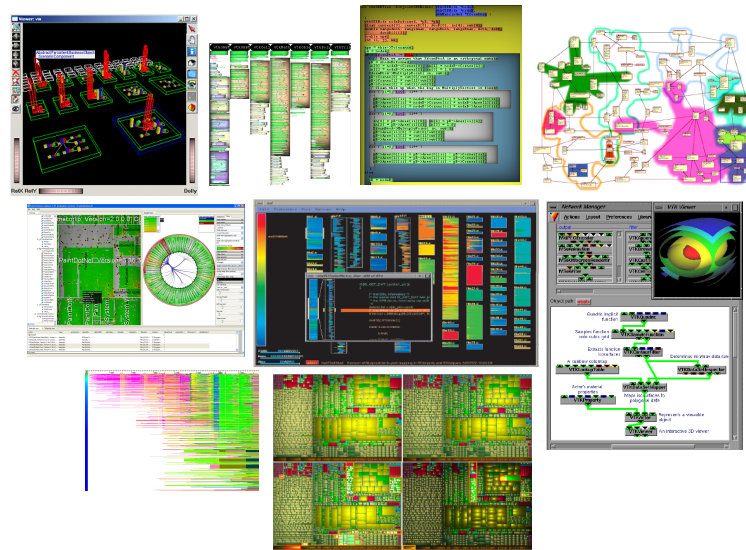


Software Visual Analytics for Maintenance

Example Solutions



prof. dr. Alexandru (Alex) Telea

Department of Mathematics and Computer Science
University of Groningen, the Netherlands

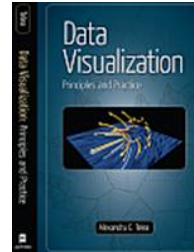
Introduction



www.cs.rug.nl/~alex



www.solidsourceit.com



Data Visualization: Principles and Practice A. K. Peters, 2008

Professor of Computer Science (Multiscale Visual Analytics),
University of Groningen, the Netherlands

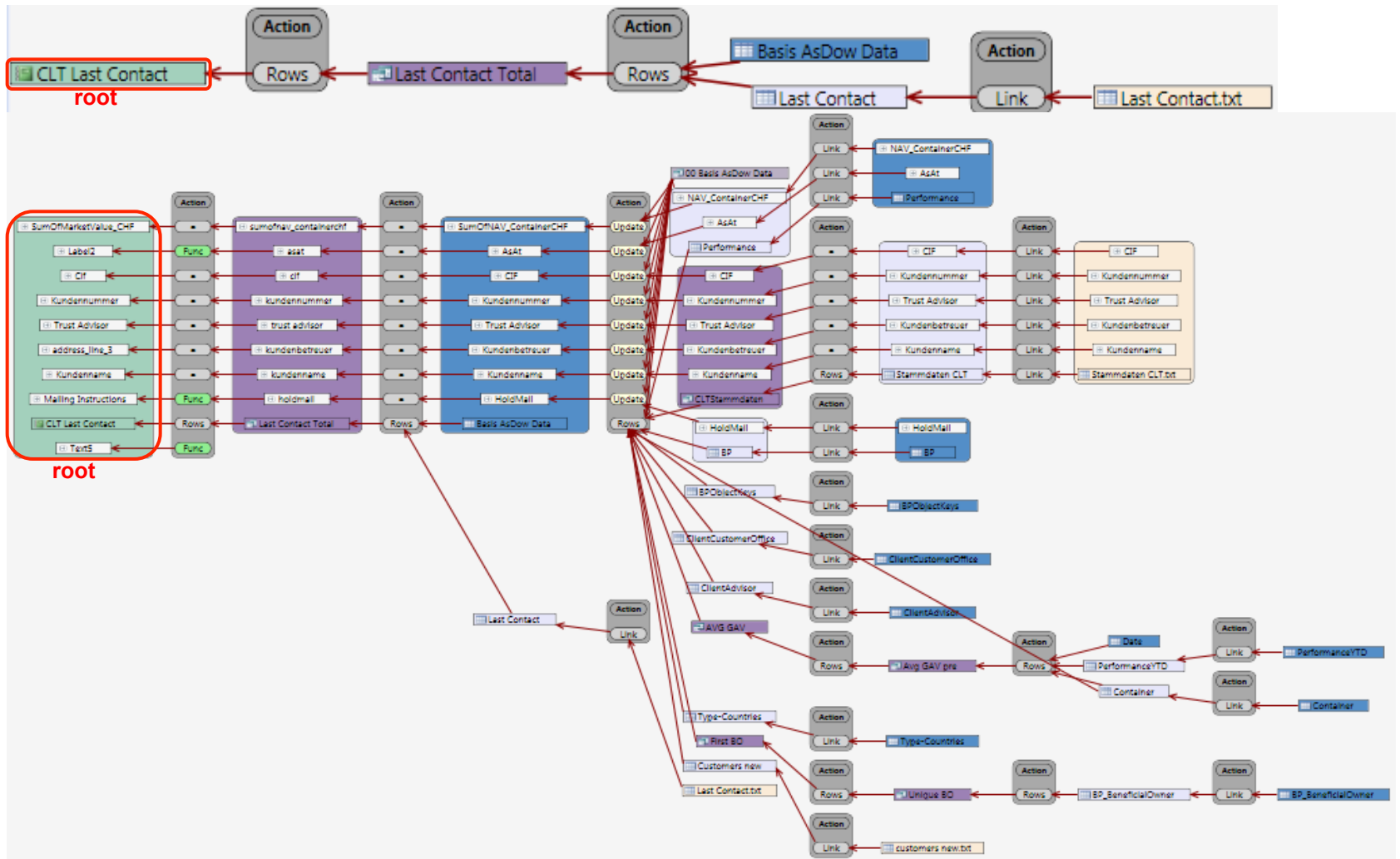
My PhD students...



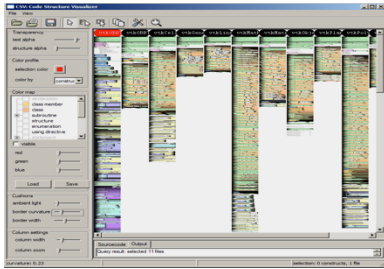
My MSc students...

Avdo Hanjalic, Tijmen Klein, Johan v/d Geest, Mark Ettema, Daniel Kok, Karsten Westra, Yuri Meiburg, Hessel Hoogendorp, Liewe Kwakman, Madalina Florean, Bertjan Broeksema, Mark Stoetzer, Sergio Moreta, Kees van Koten, Frans Boerboom, Arjan Janssen, Freek Nossin, Matthijs van Eede, Martijn van Dortmund, Maurice Termeer, Iwan Vosloo, Gerard Lommerse, Dennie Reniers, Milan Pastrnak, ...

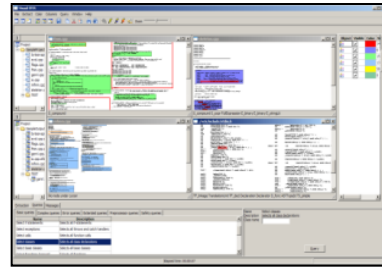
Software Visualization?



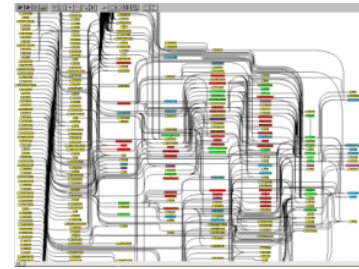
Software Visualization!



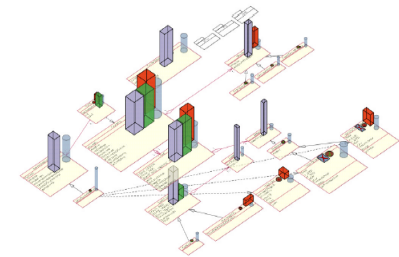
source code



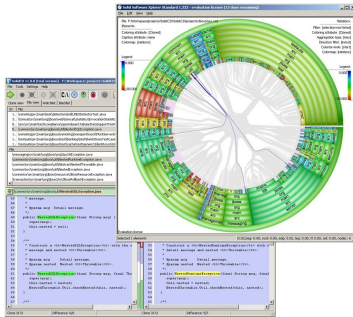
code quality



code dependencies



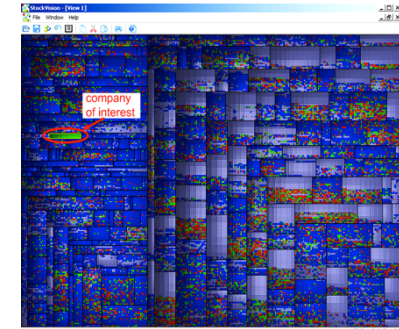
design and metrics



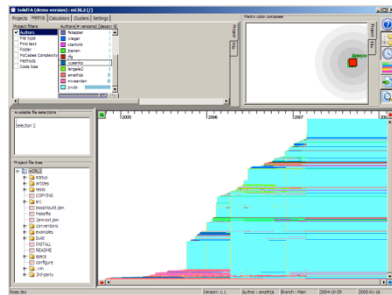
text duplication

How should we deal with **scale**?

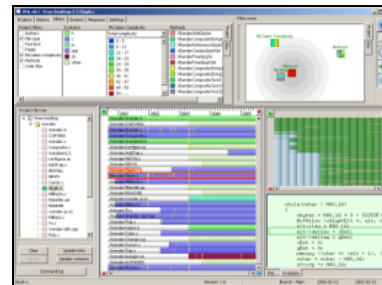
- simplified visualizations?
- continuous simplification?
- what to simplify exactly?
- reinvent wheel for each app?



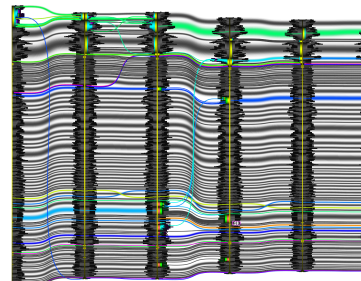
program dynamics



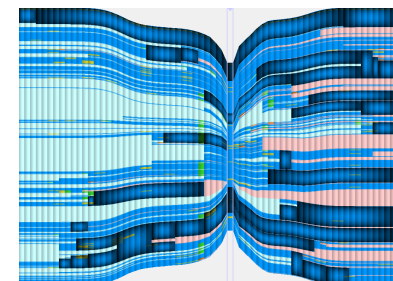
code repositories



evolution metrics



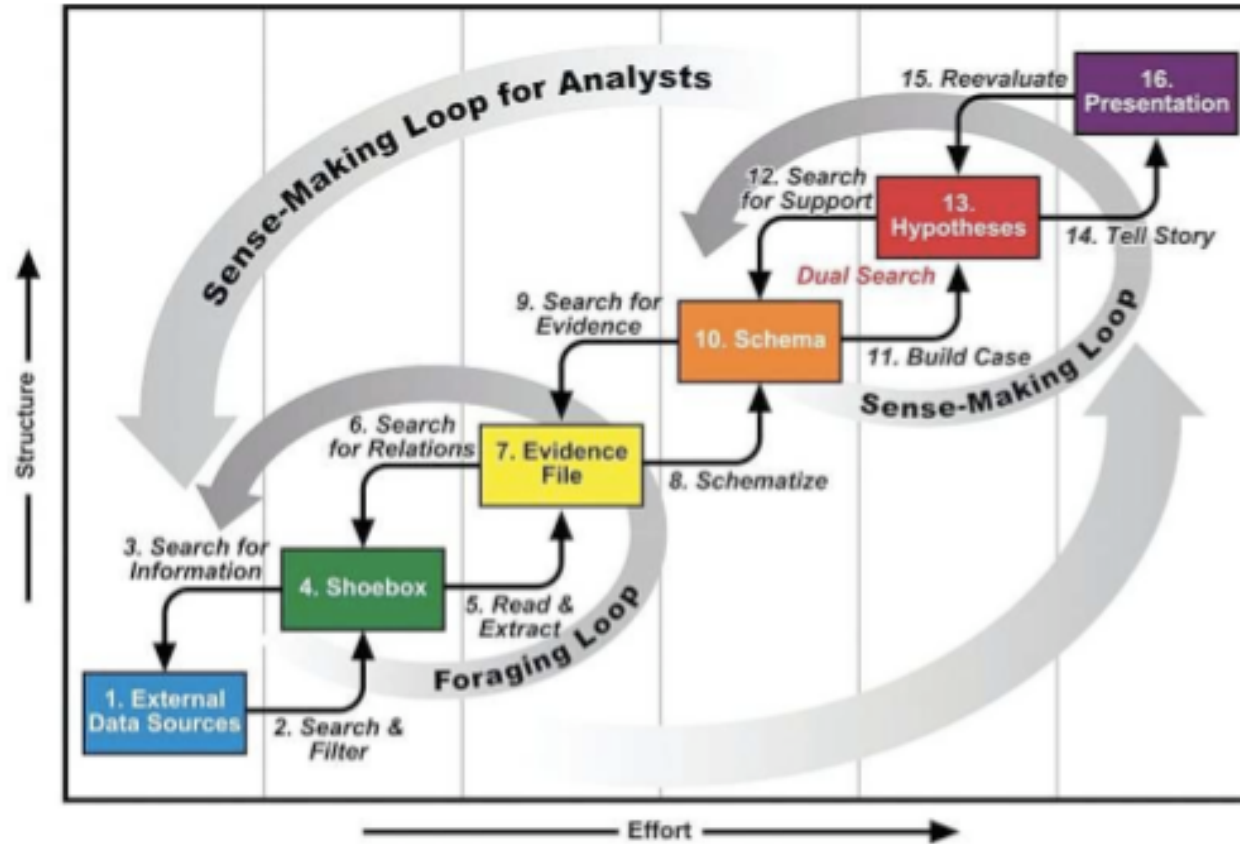
structure evolution



team analysis

Software Visual Analytics – Process View

“The science of analytical reasoning facilitated by interactive visual interfaces”



The Sensemaking Loop

- going from **raw data** to **meaning** (semantics) to **insight** to **decisions**
- data → hypothesis → (in)validation → conclusions → presentation
- put simply: **combine analysis and visualization**

Software Visualization

Definition:

- “The static or animated 2D or 3D visual representation of information about software systems based on their structure, history, or behavior in order to **help software engineering tasks**” [Diehl, 2006]

Surveys:

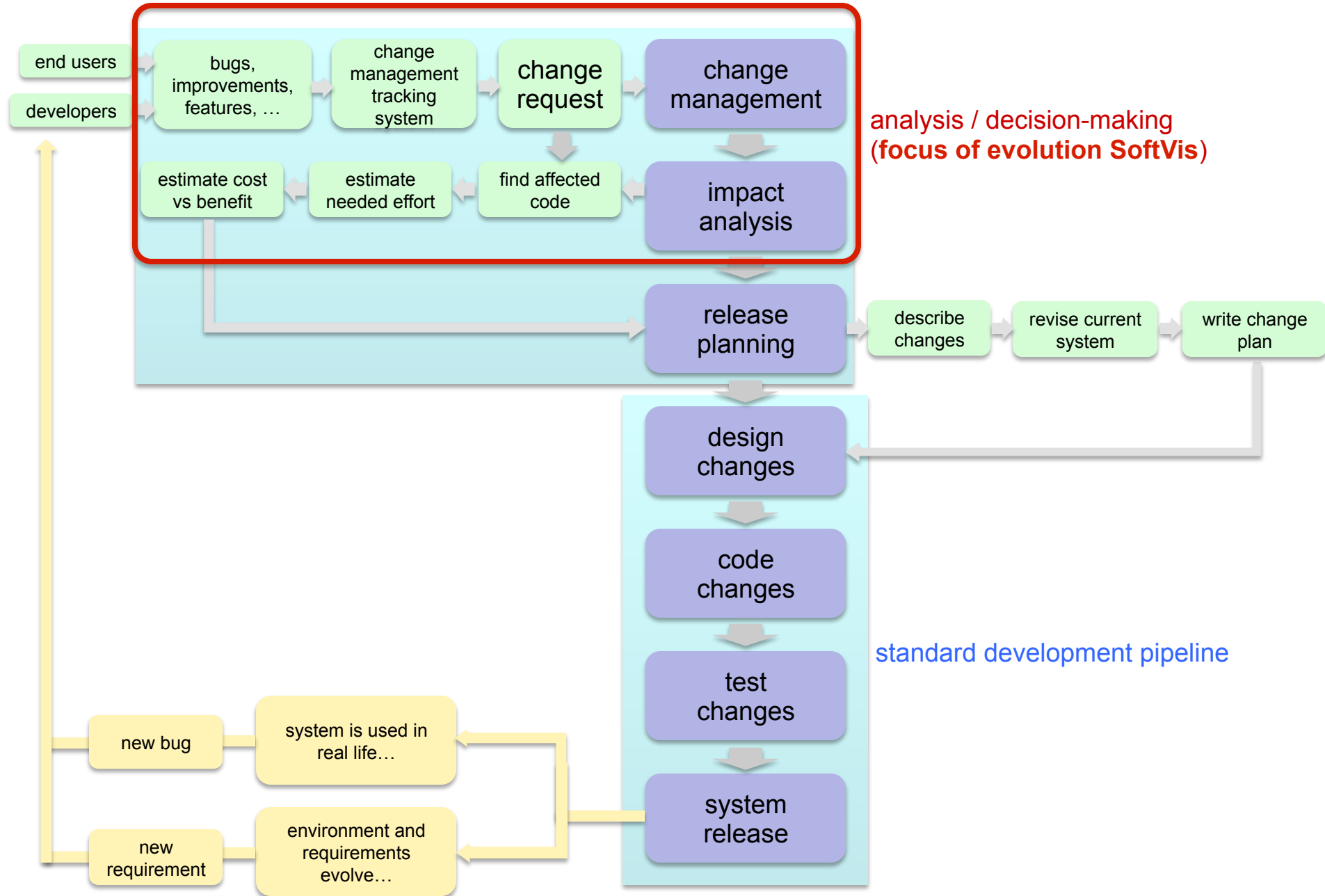
- IT industry: 457 billion \$ (2013), **50% larger than in 2008** [www.infoedge.com]
- comparison: total US health care spending 2.5 trillion \$ (2009) [www.usatoday.com/news/health]
- **80% of development costs** spent on maintenance [Standish' 84, Corbi' 99]
- **50% of this** is spent for **understanding** the software!

Practice:

- **40% engineers find SoftVis indispensable**, 42% find it not critical [Koschke '02]

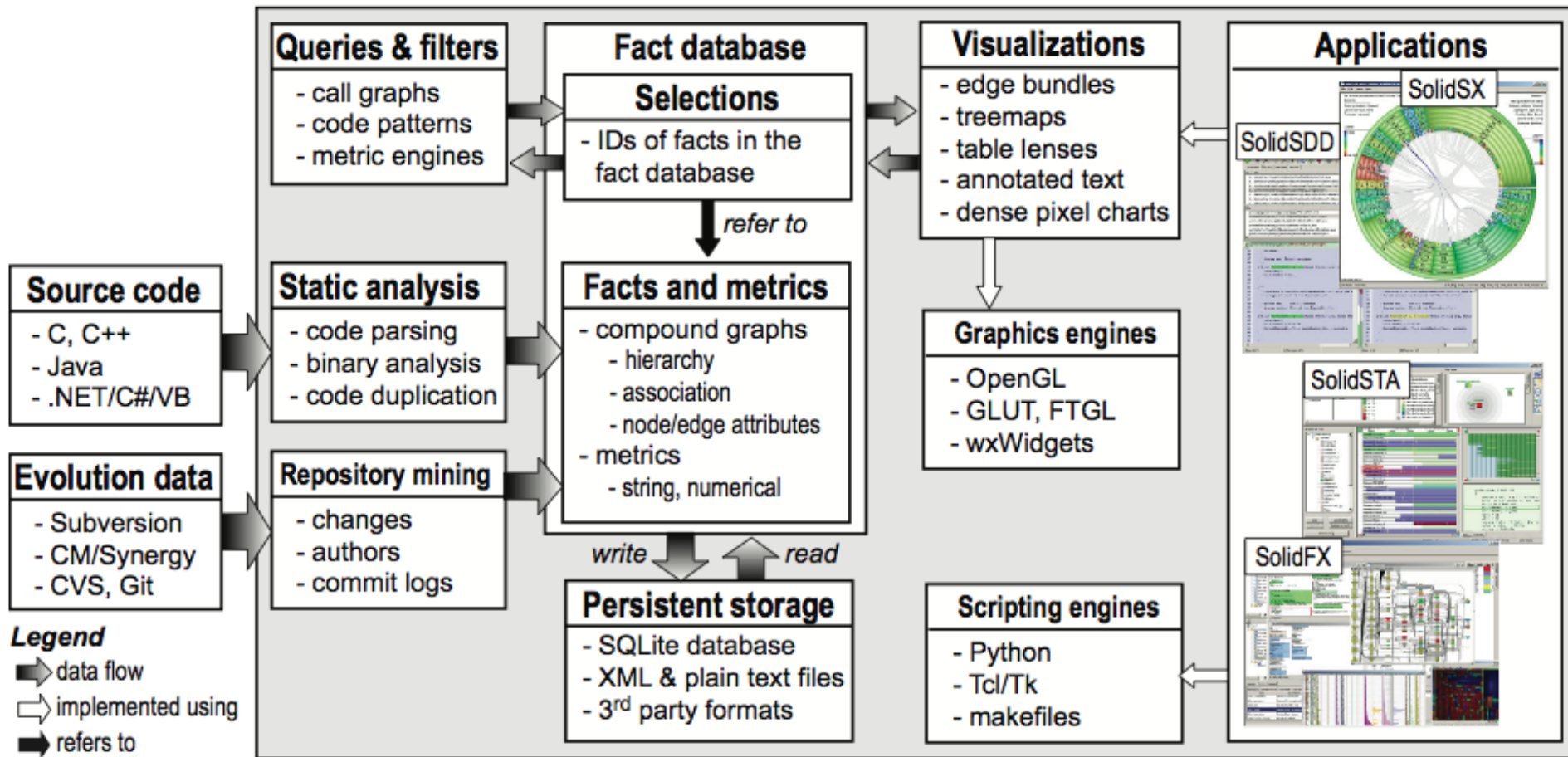
Goals: reduce cost/time, increase quality and productivity!

Visual Analytics in Software Maintenance



Techniques

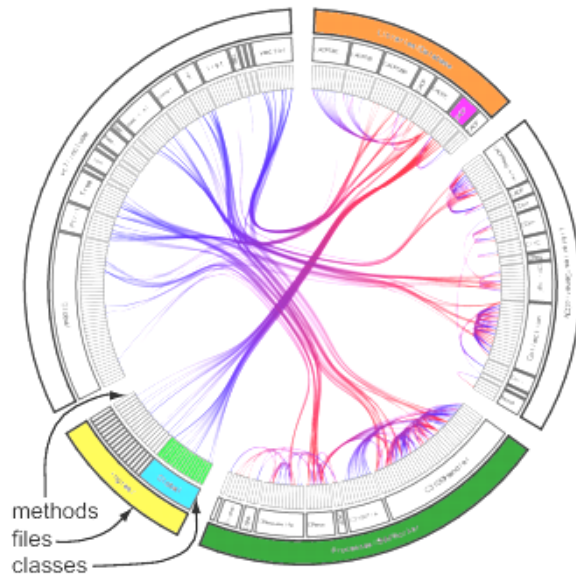
Software Visual Analytics – Technical View



Many types of **data** and **questions** → many types of **visualizations**

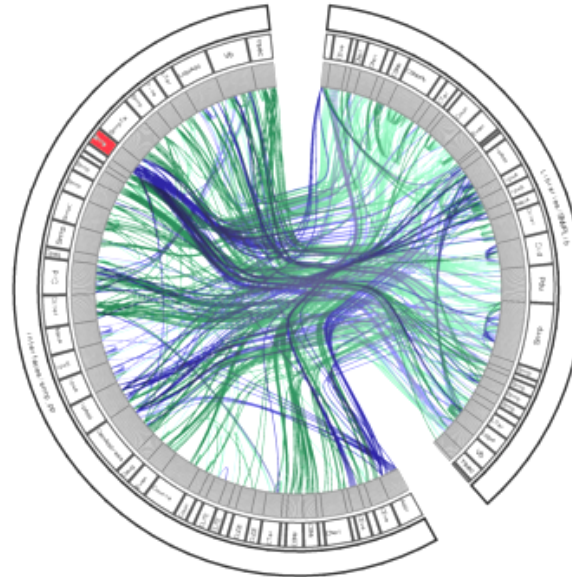
1. Assessing system modularity

Modular system



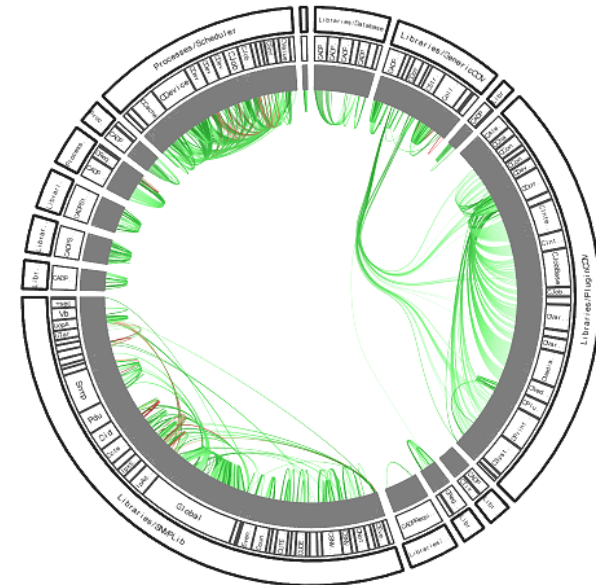
- blue = caller, red = called
- all functions in the yellow file call the purple class
- green file has many self-calls

Monolithic system



- blue = virtual, green = static functions
- red class has many virtual calls (possible interface class)

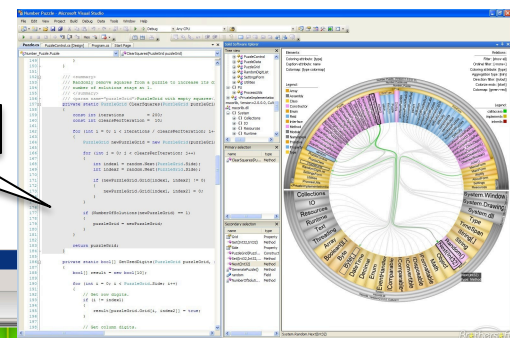
Decoupled system



- many intra-module calls
- few inter-module calls
- typical for library software

2. Structure, dependencies, metrics

SolidSX analytics tool (www.solidsourceit.com)



Code view

The main interface of SolidSX is divided into several panels:

- Structure:** A tree view on the left showing the project's file and folder hierarchy.
- Treemap:** A central heatmap where the color of each cell represents a metric for a specific node or relationship.
- Node table view:** A table at the bottom left providing detailed metrics for individual nodes.
- Radial view:** A circular dependency graph at the bottom right showing the relationships between nodes.

Structure

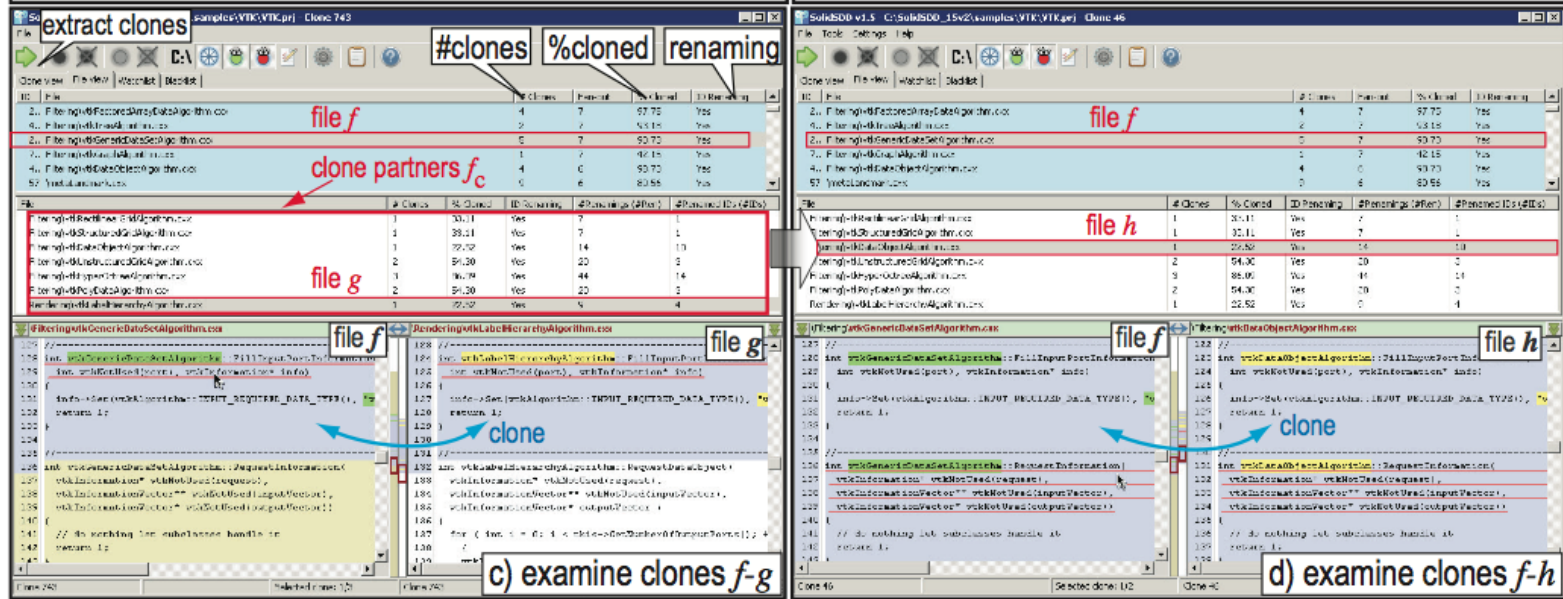
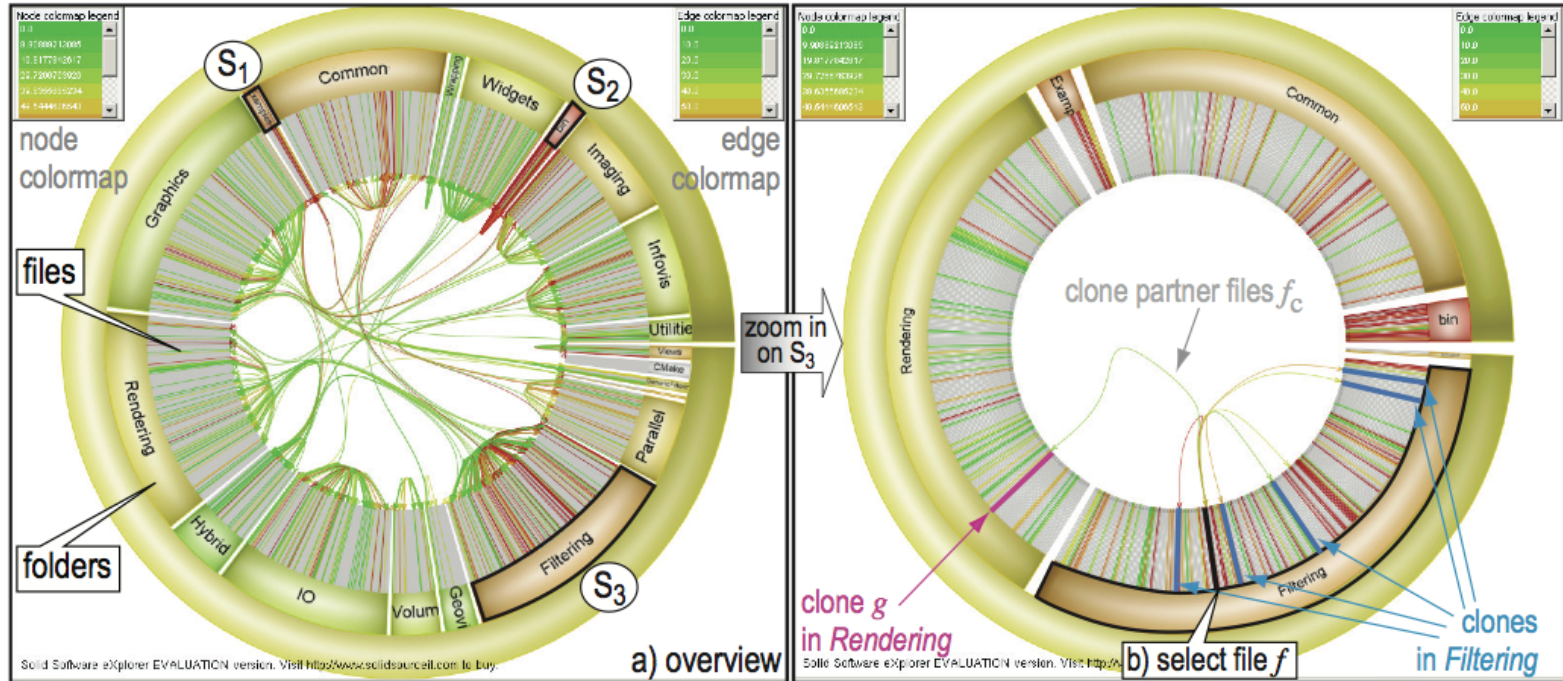
Test results

Detail metrics

Dependencies

3. Code duplication

SolidSDD tool (www.solidsourceit.com)

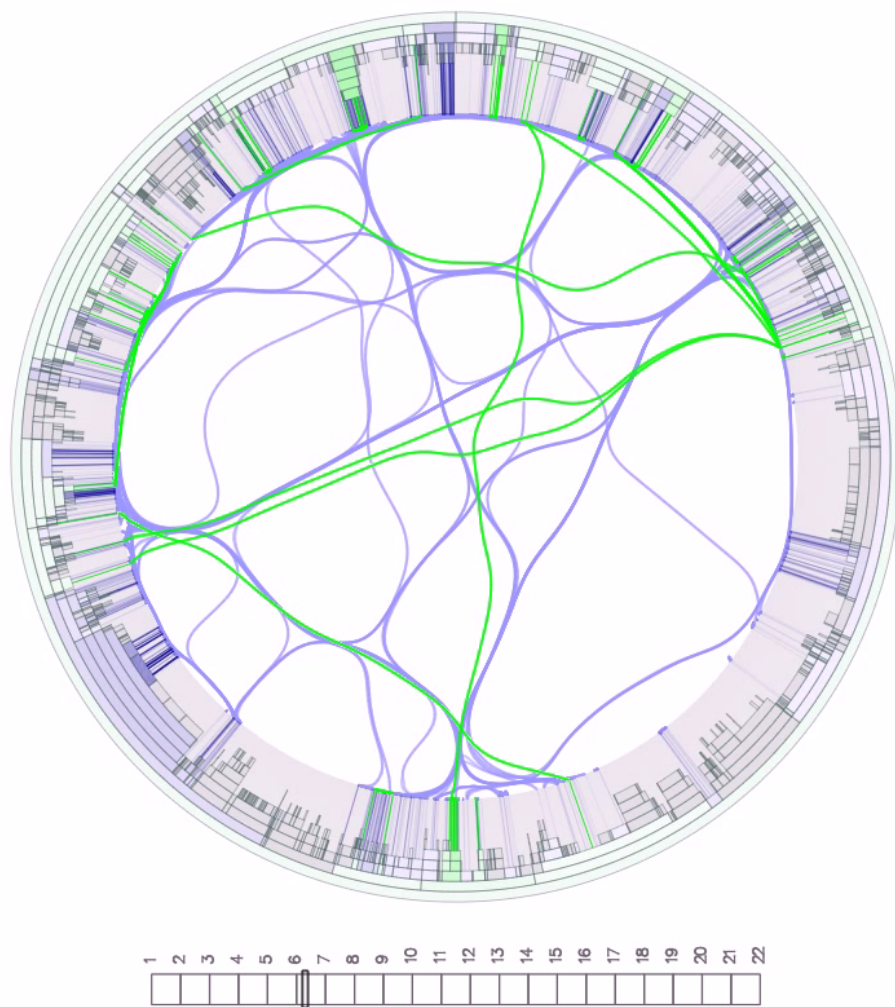


Color legend: non-cloned code clone partner not shown clone shown in both windows identifier renamed (shown as in right window)

4. Clone evolution

Questions

- how does code duplication change in time?
- which clones are added, removed, merged, or split? And why?

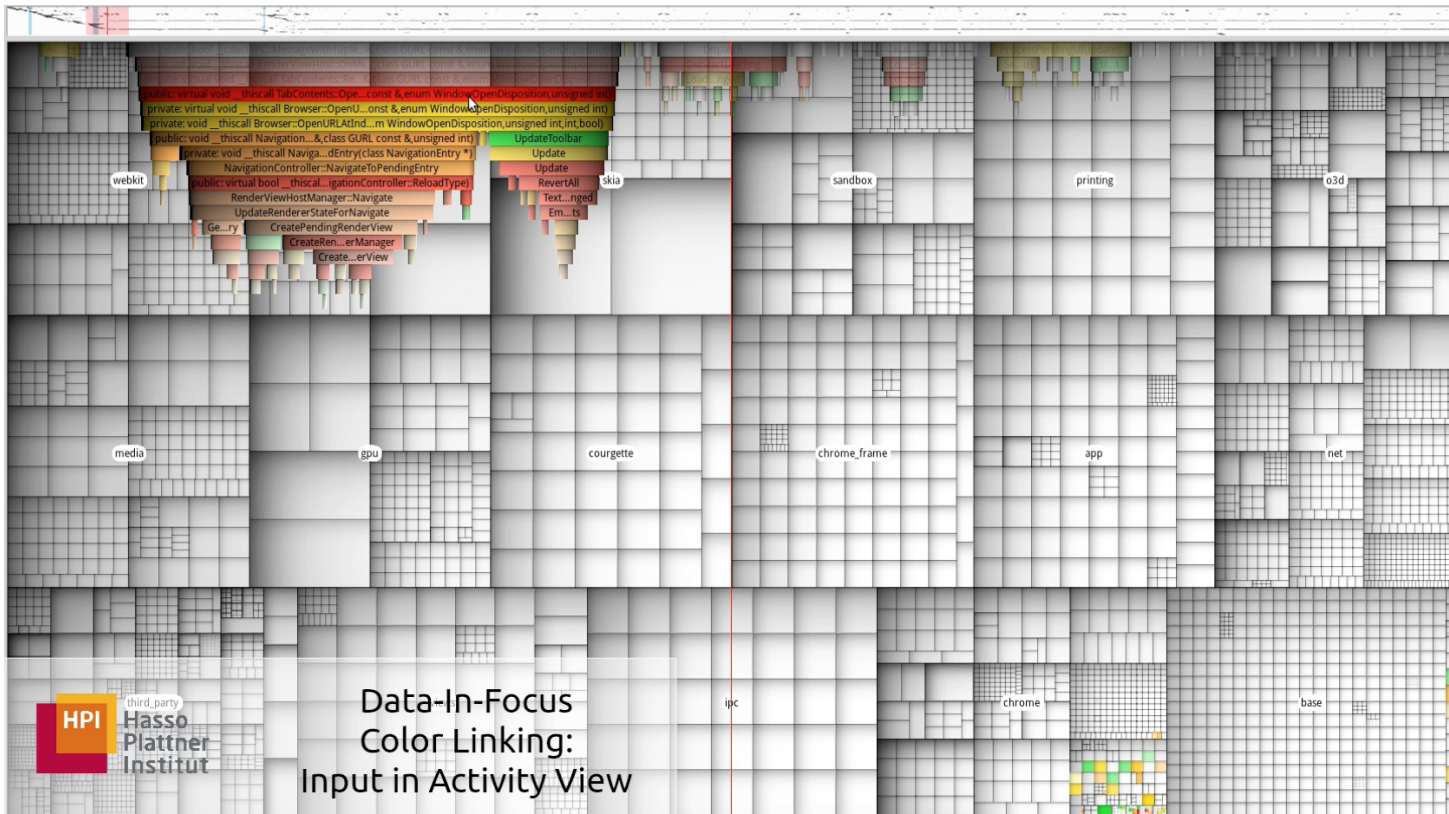


Evolution of clones in Mozilla Firefox (~55K clone relations, 3.5 MLOC C/C++)

5. Program trace and structure

Questions

- where (in the program structure) are the calls executed now?
- when (during execution) are calls to this subsystem done?



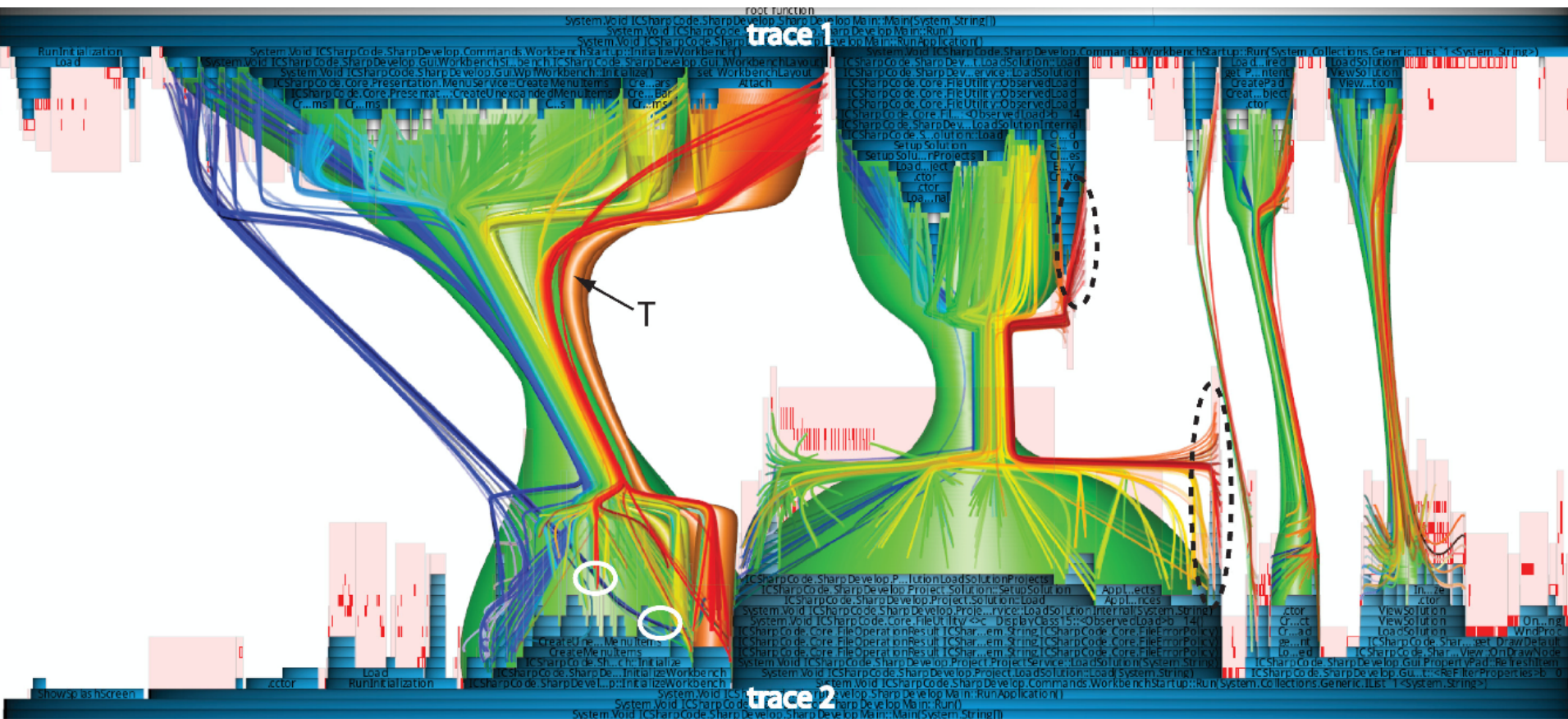
Code: Chrome browser (2.7 MLOC C/C++, 8900 files+folders)

Trace: 9000 calls to 914 functions

6. Comparing program traces

Questions

- given 2 traces, where are similar and where are different call-blocks?
- how to spot differences in call moment, duration, and called functions?

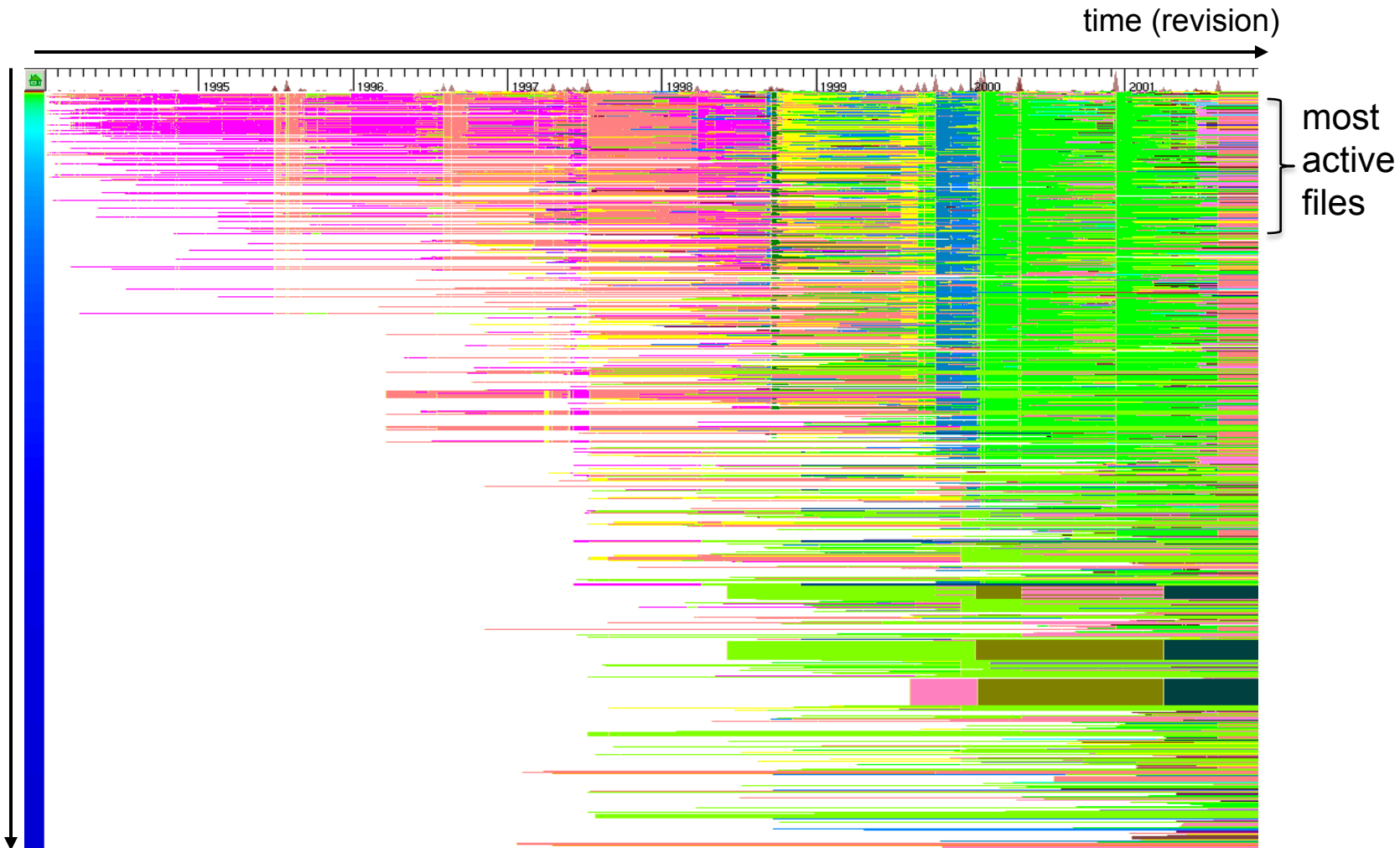


Code: 1MLOC C#, 45 developers, 8 years
Traces: 2x150K calls to 1500 functions

7. Software Evolution

Questions

- how to correlate **metrics** over large software repositories (>10K files, >100K commits?)
- how to detect **trends** to predict the future (cost, effort, risk)?

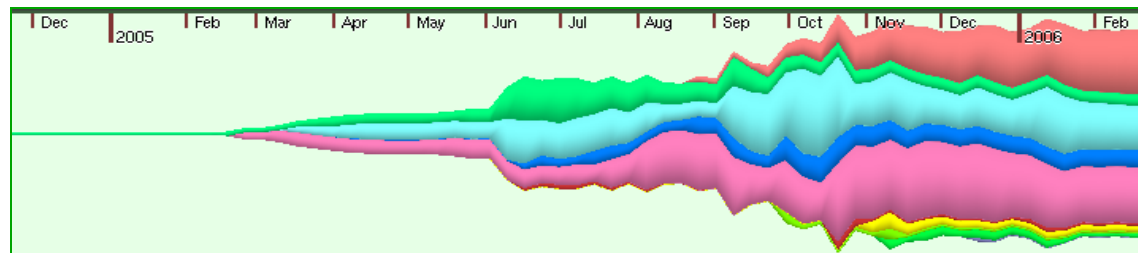


Analyzing developer effort

Show aggregated **developer impact** (#files modified by each developer) over time

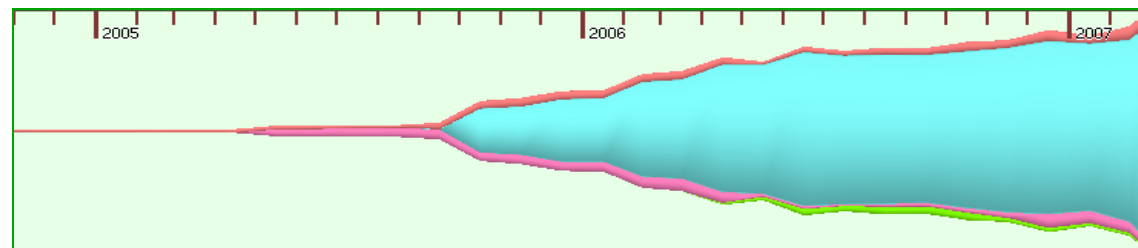
Project A (open-source)

- software grows in time
- impact: balanced over most developers



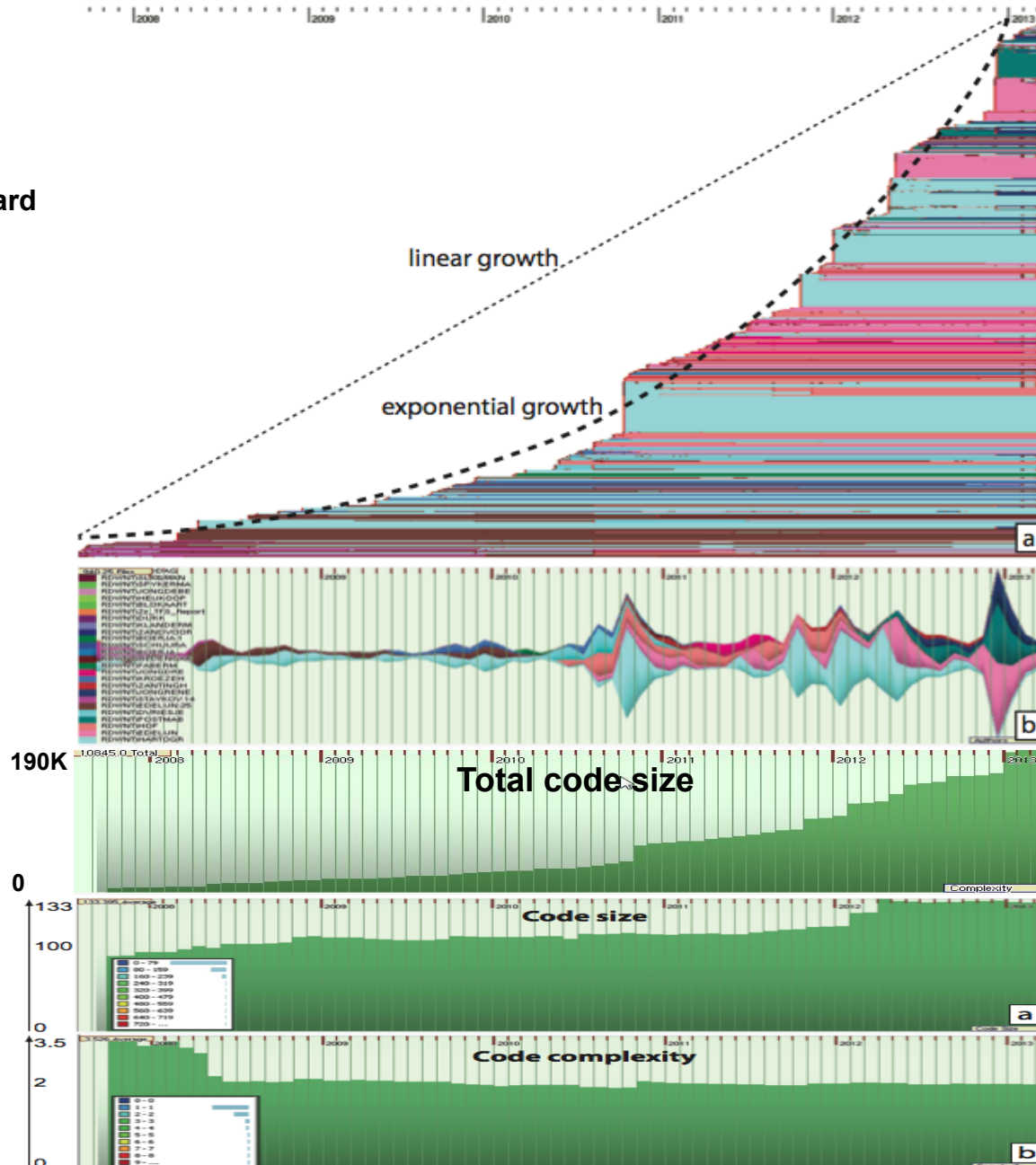
Project B (commercial)

- software grows in time at about the same rate
- but one developer owns most of the code
- what if this person **leaves** the team?!



Correlating quality metrics

- C# code base
- 4 years, 190 KLOC
- Permanent quality monitoring dashboard solution

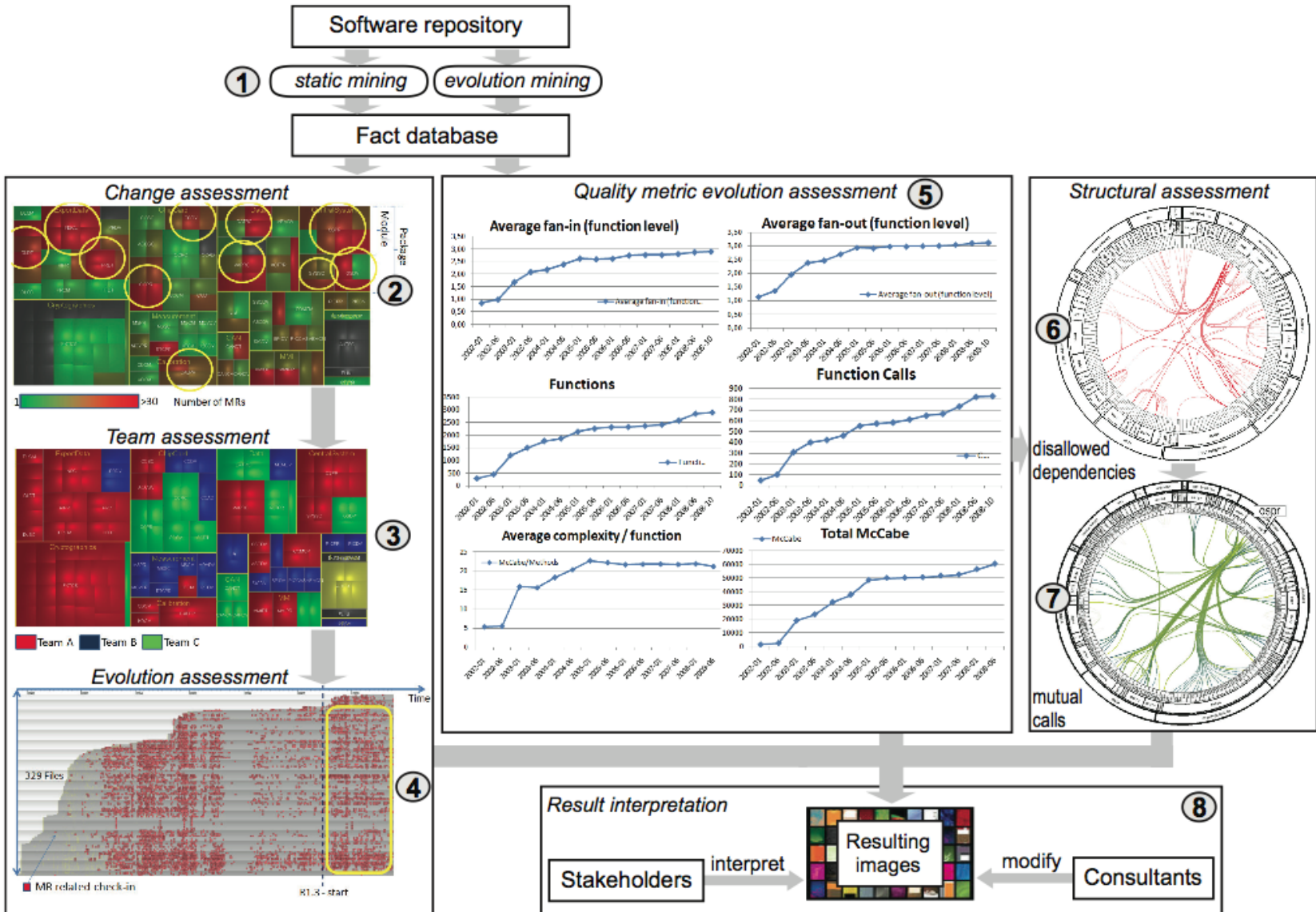


Applications

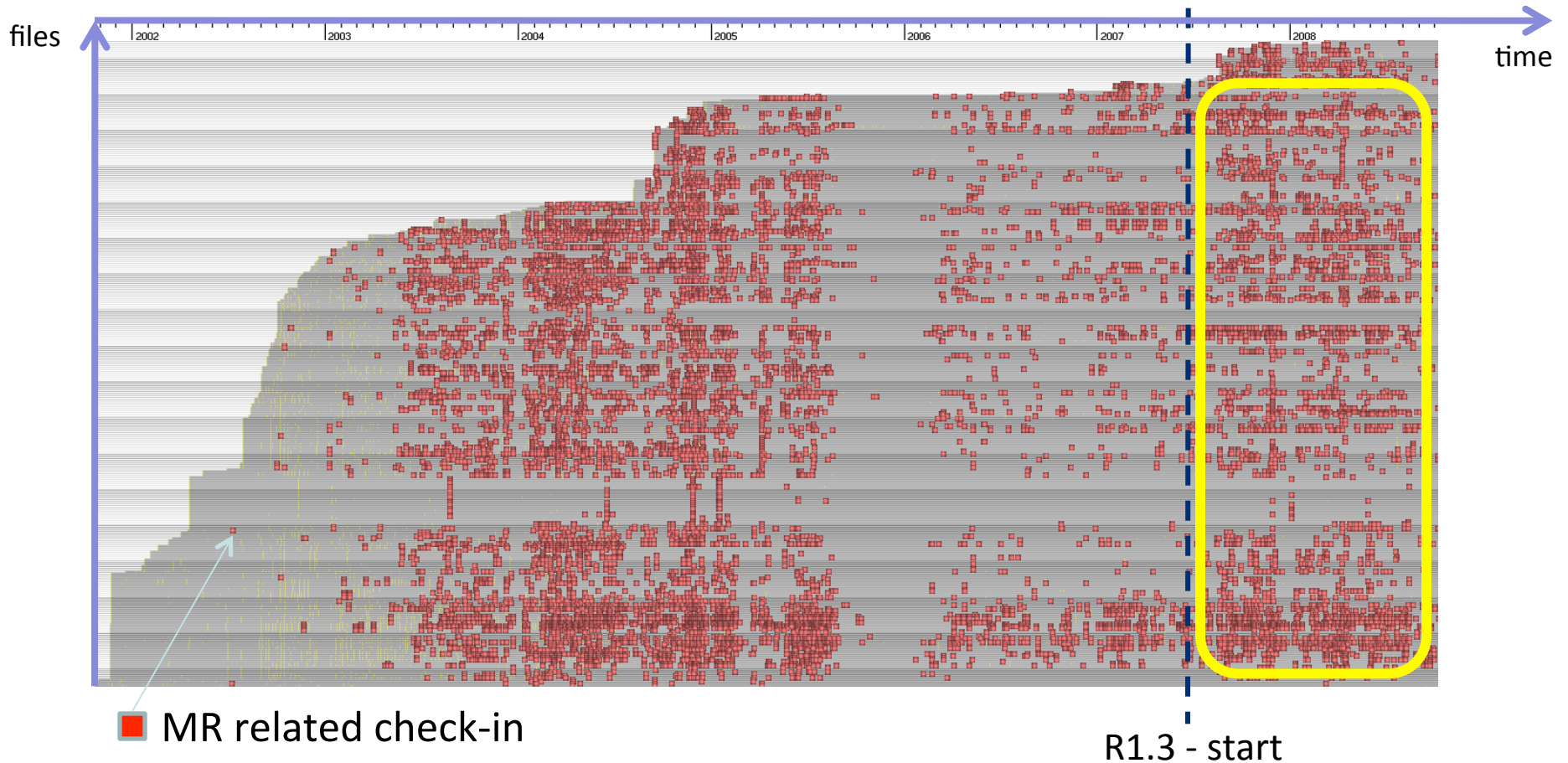
Application: Post-Mortem Assessment

Questions

- automotive project: 8 years, 3.5 MLOC embedded C, 15 releases, 60 developers
- project failed to deliver. Why?

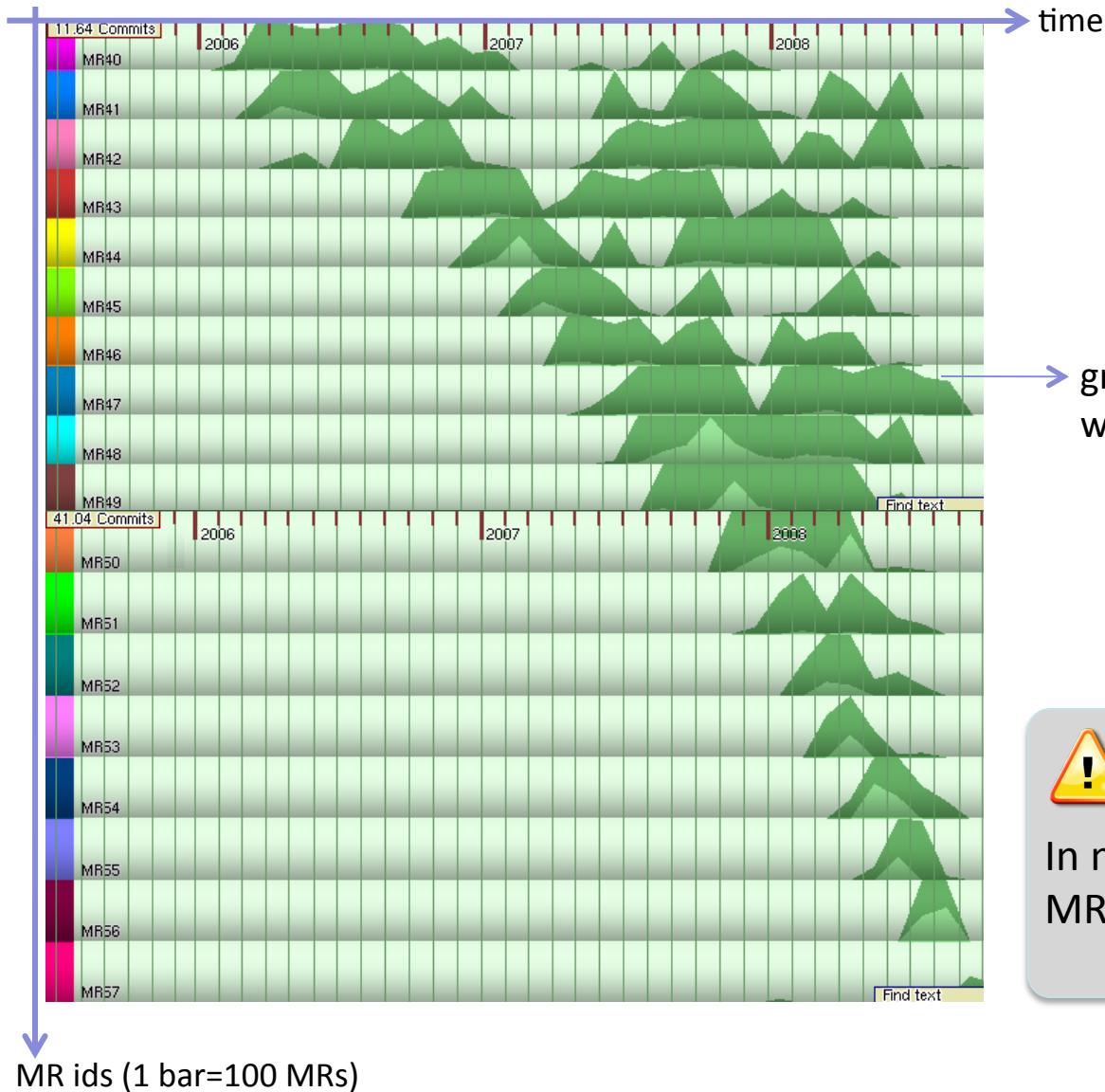


Analysis 1: Modification Request (MR) Lifetime



Little increase in the file curve – most activity in *old* files suggests too long **maintenance** & **closure** of requirements

Requirements: MR Duration

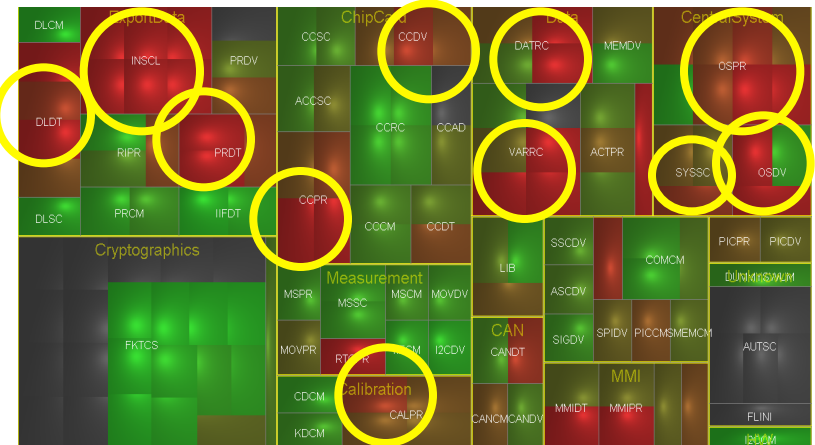


In mid 2008, activity related to MRs from 2006 *still takes place*

Analysis 2: Team Code Ownership



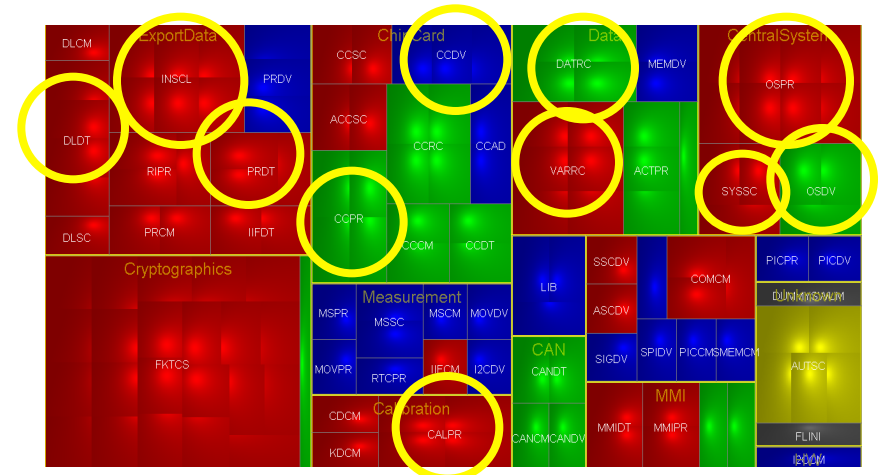
#developers
1 >8



#modification requests (MRs)
1 >30



MR closure (days)
1 >90

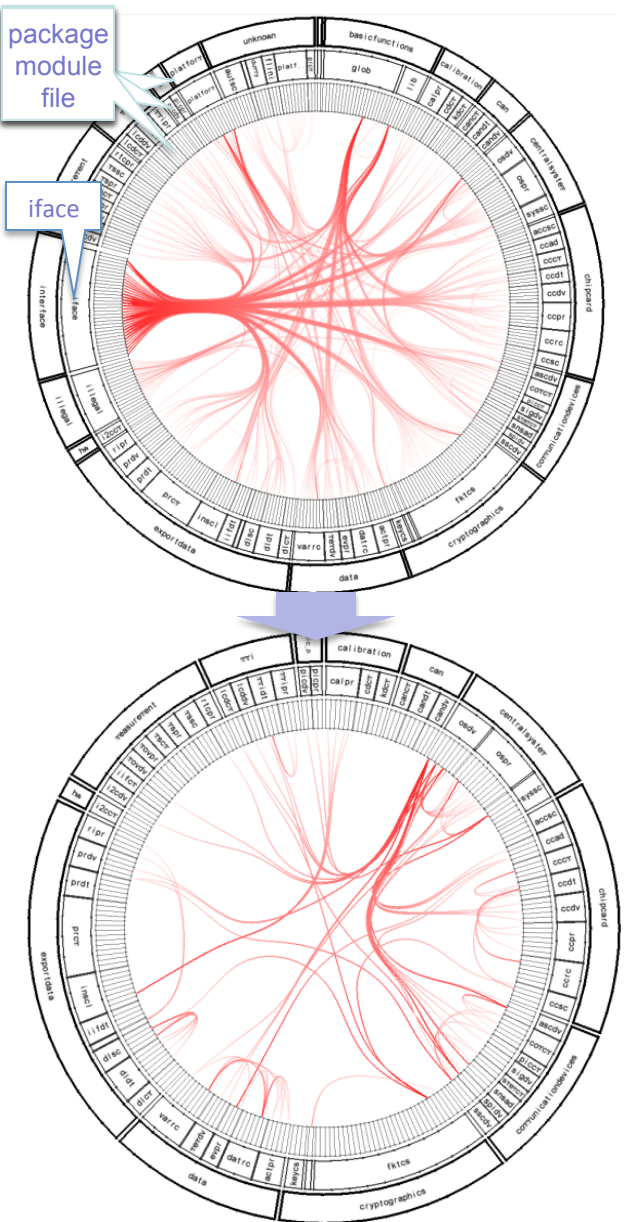


team A team B team C



Large part of software affected by long open-standing MRs
 Most of these are assigned to team A (largest team)...
 ...and this team was reported to have communication problems!

Analysis 3: Code Dependencies



uses = call, type, variable, macro, ...

is used

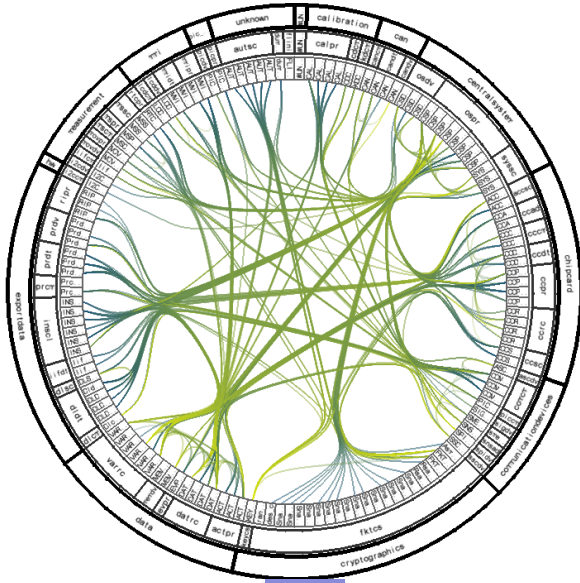
Most dependencies occur via the iface, basicfunctions and platform packages

Filter out these allowed dependencies...
...to discover *unwanted* dependencies



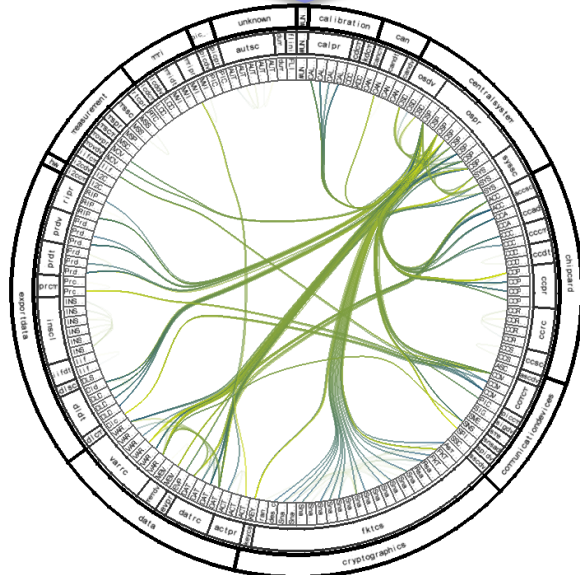
These are accesses that **bypass** established interfaces
There are several such accesses (bad)

Analysis 4: Code Call graph



High coupling at package level
This image does not tell us very much

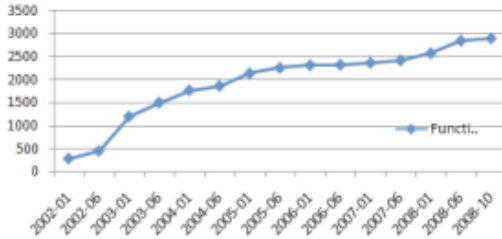
Select only modules which are *mutually call dependent*...
...to discover *layering violations*



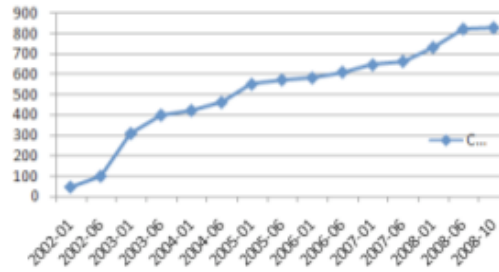
Not a **strict layering** in the system (as it should be)
Thus, the architecture is violated.

Analysis 5: Code Quality Metrics

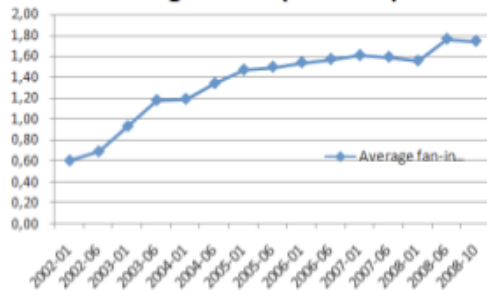
Functions



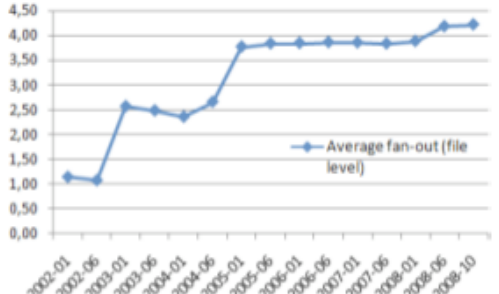
Function Calls



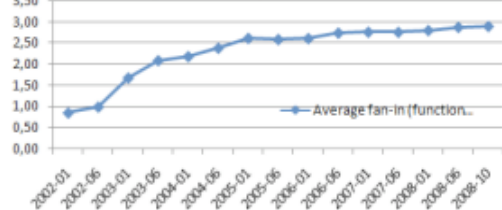
Average fan-in (file level)



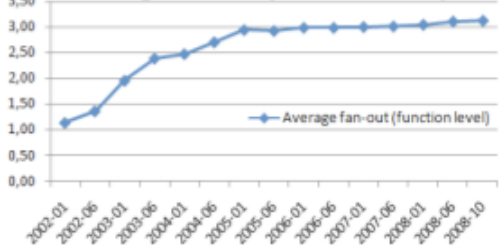
Average fan-out (file level)



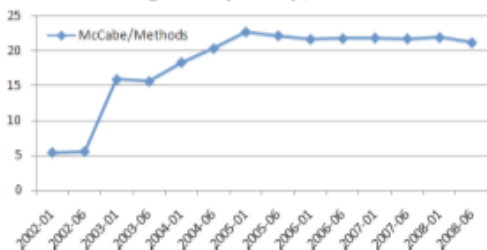
Average fan-in (function level)



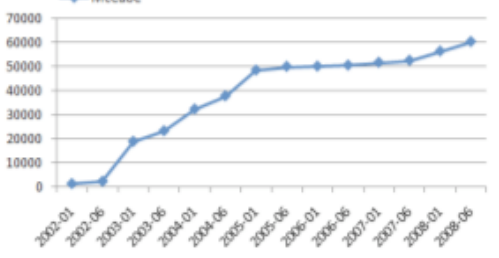
Average fan-out (function level)



Average complexity / function



Total McCabe

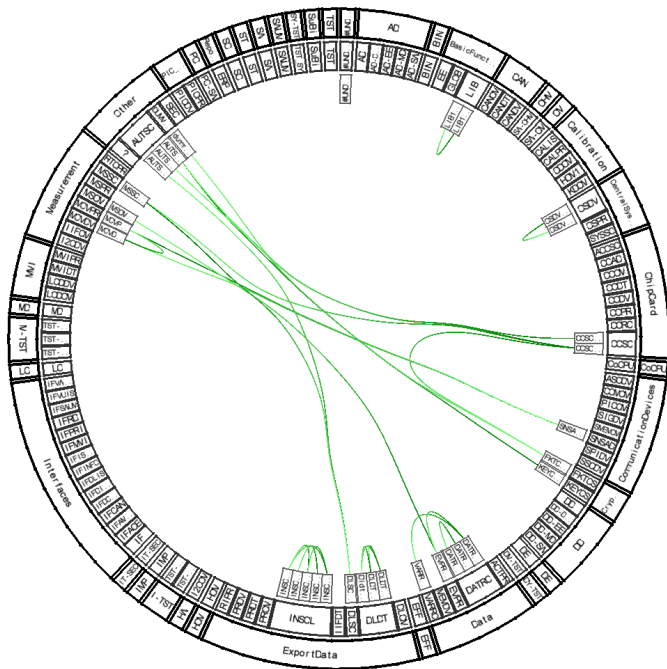


Moderate code + dependency growth
• does not explain products problems



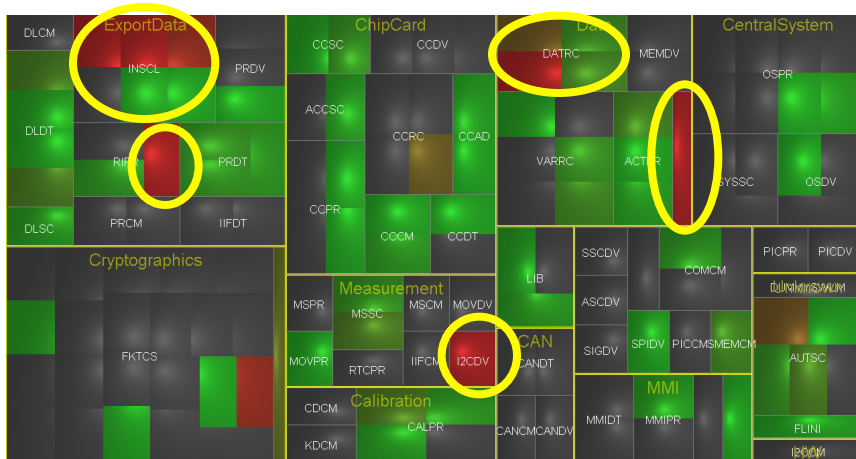
Average complexity/function > 20
Total complexity: up 20% in R1.3
• testing can be **hard!**
• **possible cause** of product's problems

Analysis 6: Code Duplication



External duplication

- show modules having similar code blocks of >25 LOC

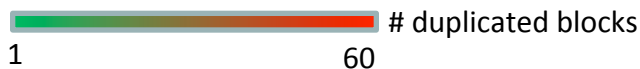


Internal duplication

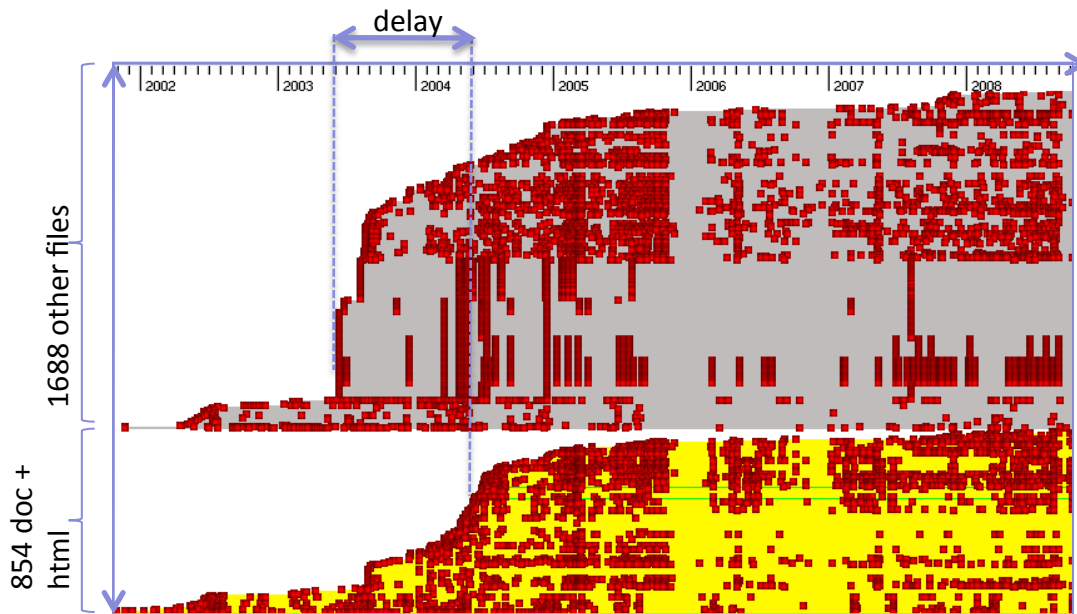
- color: #duplicated blocks within a file



Little external/internal duplication
Arguably **not a problem** for testing

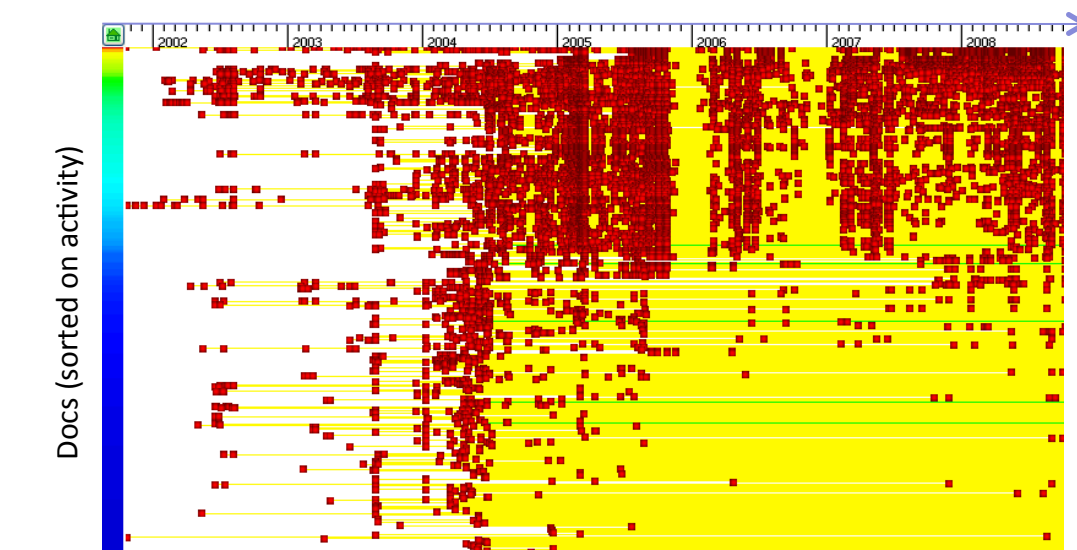


Analysis 7: Documentation



time

- **30% of files** are documentation
- updated **regularly**
- grow **in sync** with rest of code base



time

- **40% of docs** frequently updated
- rest seem to be **stale**



Code is **well documented**...
...so refactoring likely doable
Start from up-to-date docs

Application: Database reverse engineering

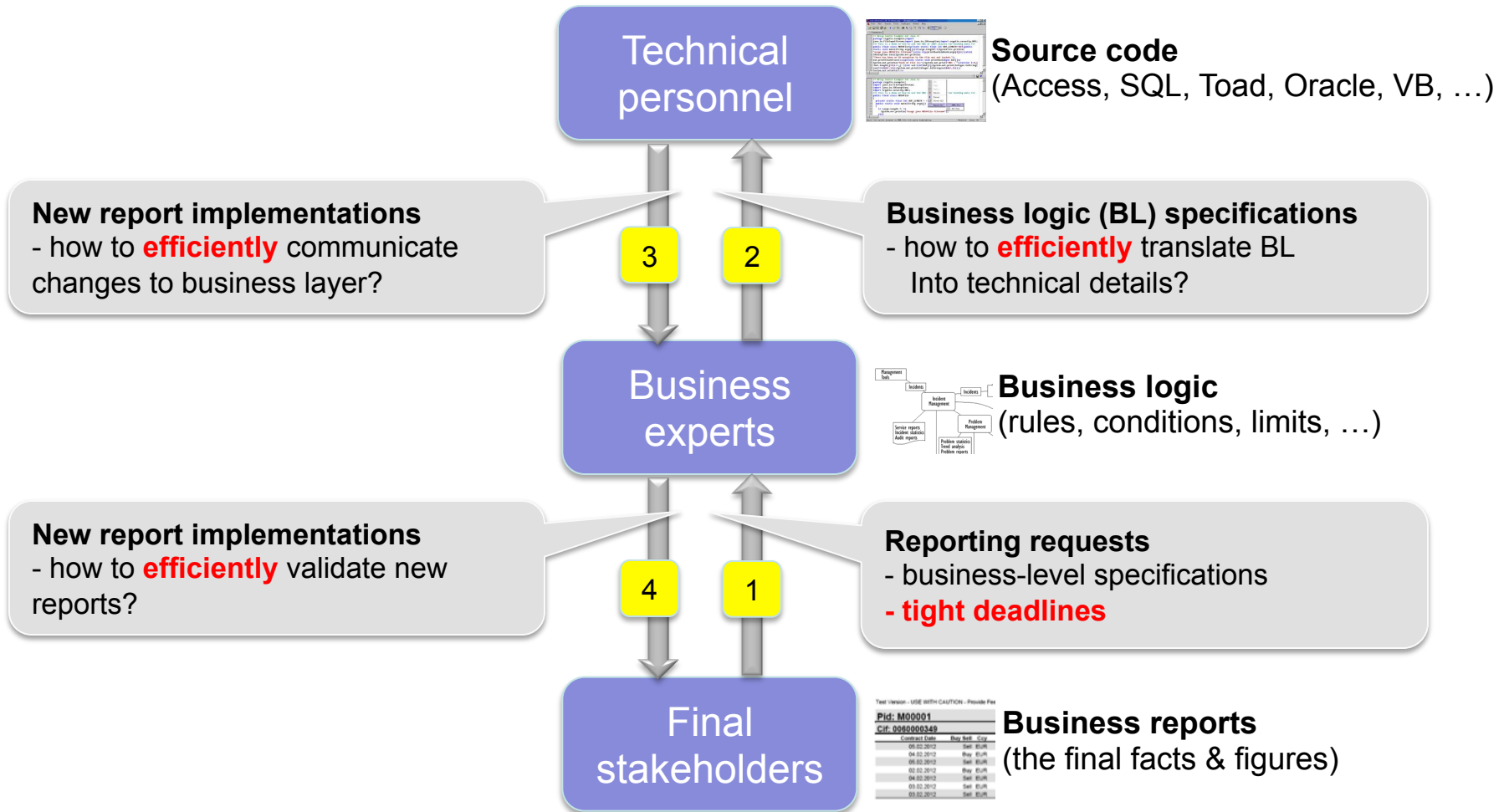
Context

- clients: top-3 Swiss bank
- product: reporting system (2004-2012)
 - Oracle/SQL/MS Access databases
 - ~5000 tables, 60000 fields,
 - mix of TS-SQL, Visual Basic, MS Access
 - code needs **24-hour uptime**
- system **was unmaintainable**, at end

Questions

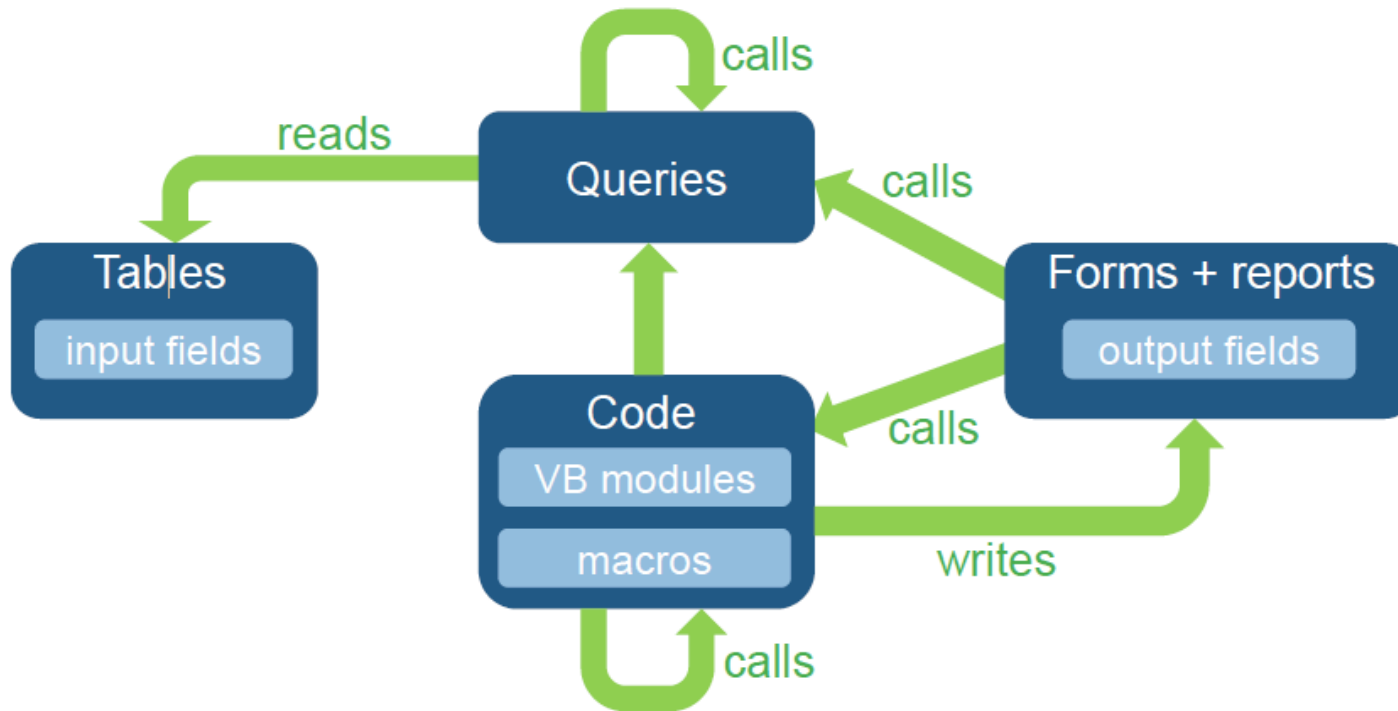
- how can we understand the **business logic**?
- how can we refactor the database design for better **maintenance**?

Stakeholders



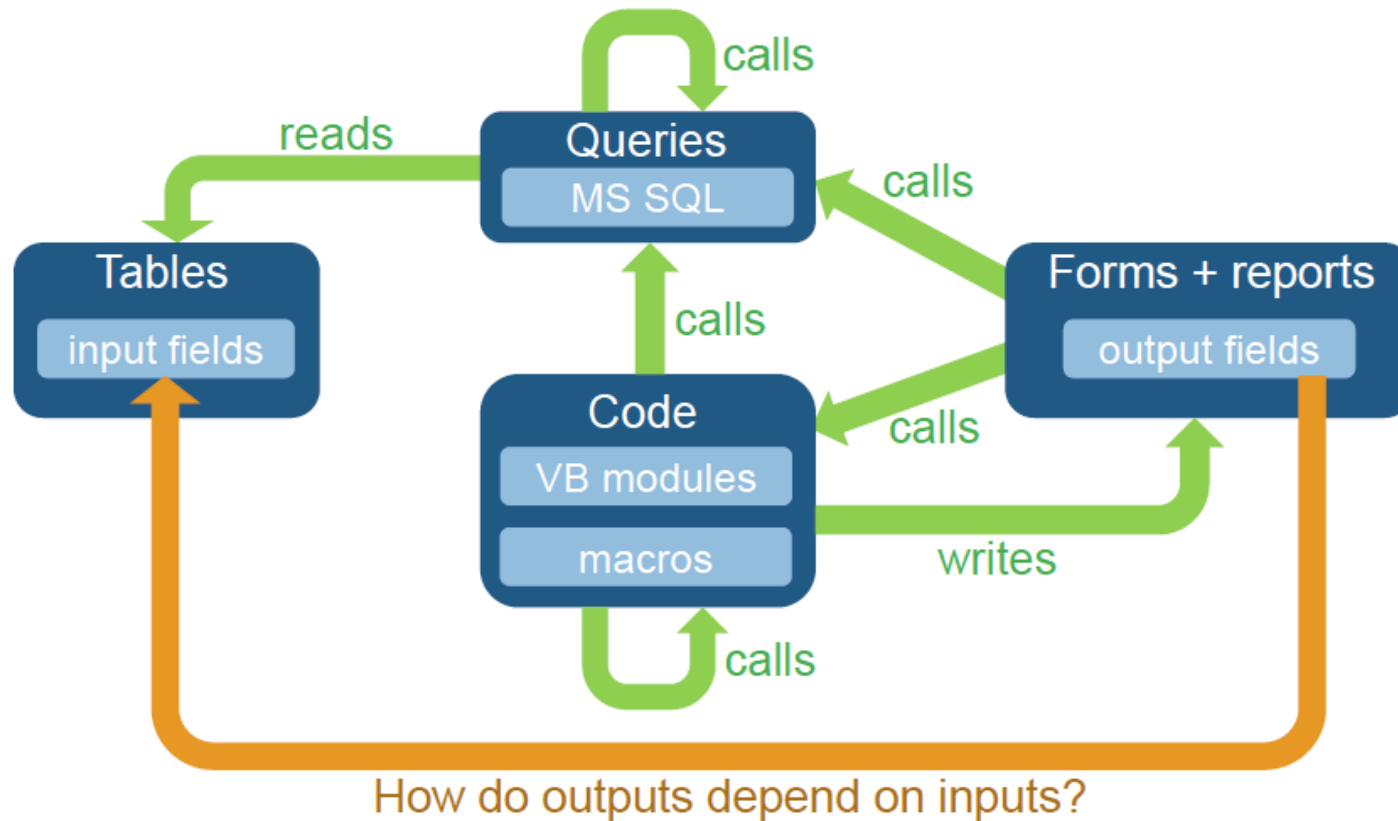
Problem Modeling

Data and control dependencies



Problem Modeling

Data and control dependencies



- one output to many input, one input to many output relations
- hard to find all relations purely statically – dependence on control flow (execution paths)

Problem Dimension

Entity	Measurements (AsDow)
databases processed	39
tables	1'258; 11'235 input fields
reports	66; 1'396 fields; ~21 fields/report; ~8 input/output relations (min.)
queries	3'617
VB code	52 modules; 271 functions; 24 KLOC
macros	22

Entity	Measurements (Reports)
databases processed	58
tables	1'977; 24'474 input fields
reports	351; 5'946 fields; ~17 fields/report; ~4 input/output relations (min.)
queries	2'010
VB code	69 modules; 186 functions; 17 KLOC
macros	110

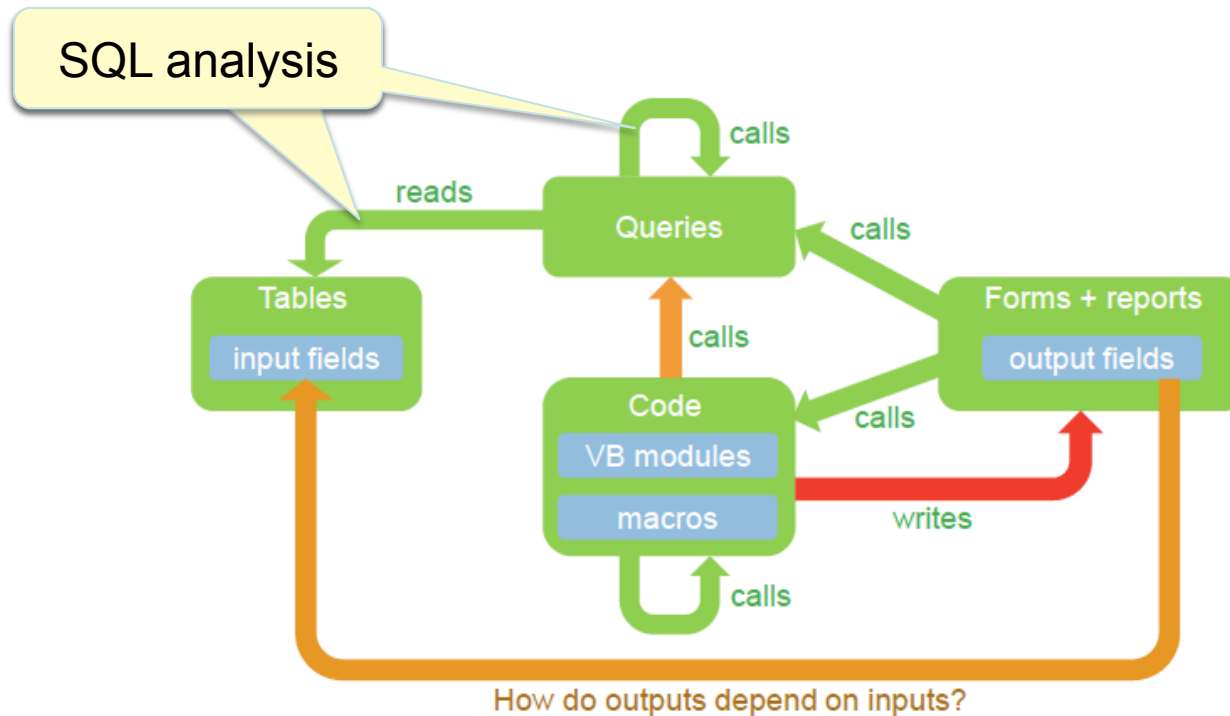
Input parameters

- 35709 input fields
- 5..10 processing steps per path
- assume input fan-out of 4 – very conservative
- total: $4 \times 35709 = 142836$ output-to-input field paths

Manual effort needed for impact analysis

- 10 mins/path needed – very conservative
- result: 23806 hours ~ **15 person years!**

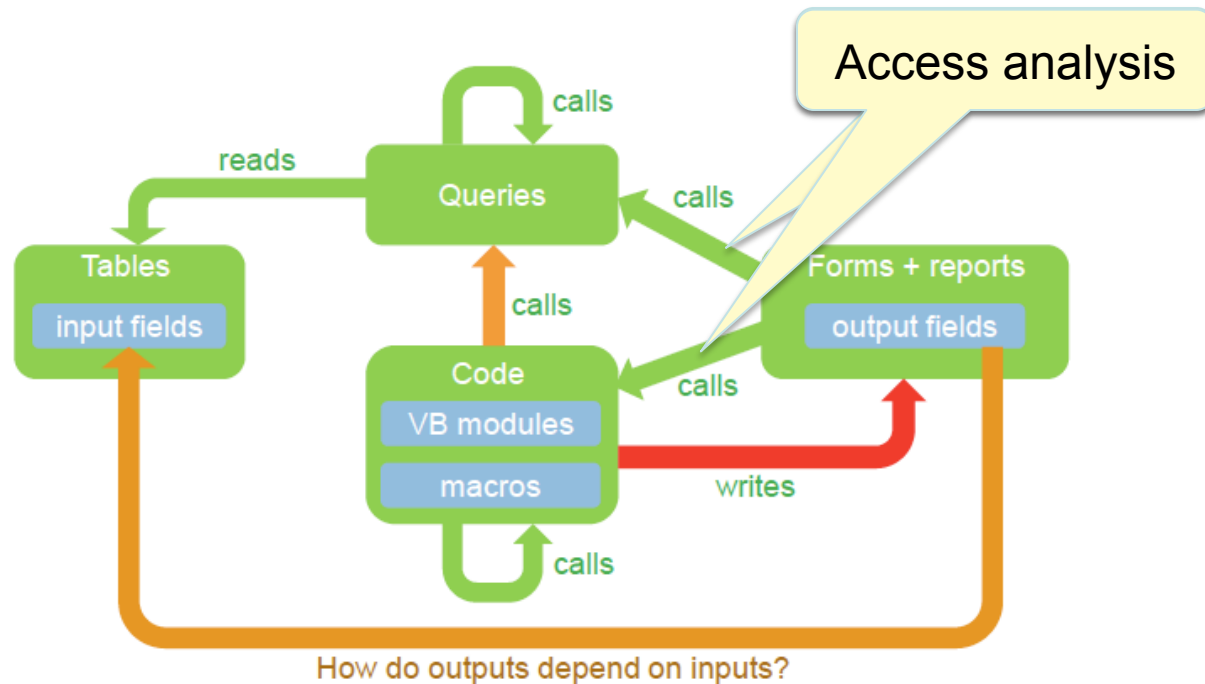
Solution - step 1



SQL analysis

- extract all SQL code (TS-SQL, Oracle, MS Access scripts)
- perform a syntax analysis on SQL (parsing)
- find all read table names and fields (columns)
 - use control flow analysis of SQL

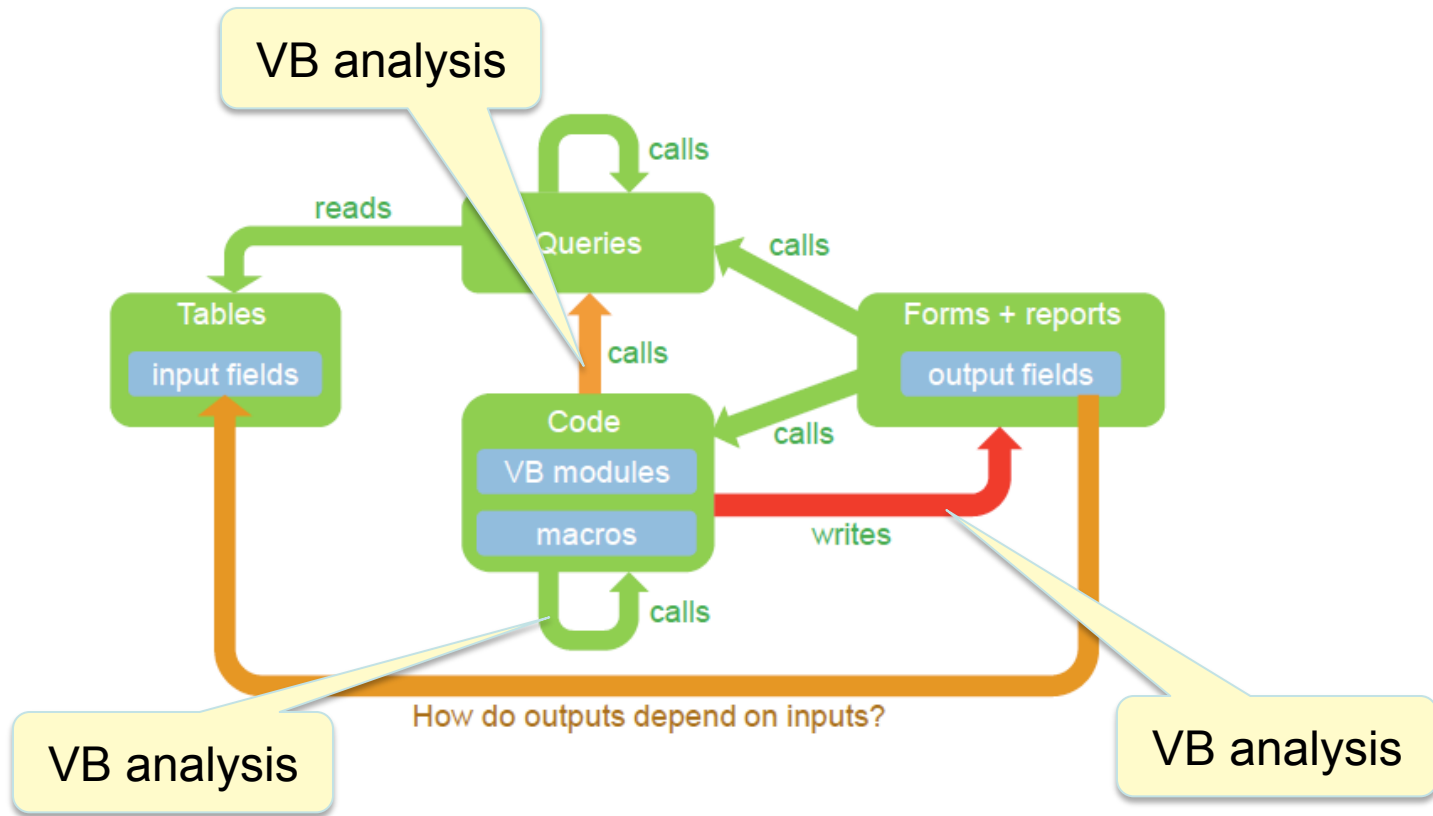
Solution - step 2



MS Access analysis

- extract all MS Access form (report) definitions
- identify fields of interest
- find SQL or VB code that these fields call (if any)

Solution - step 3



VB analysis

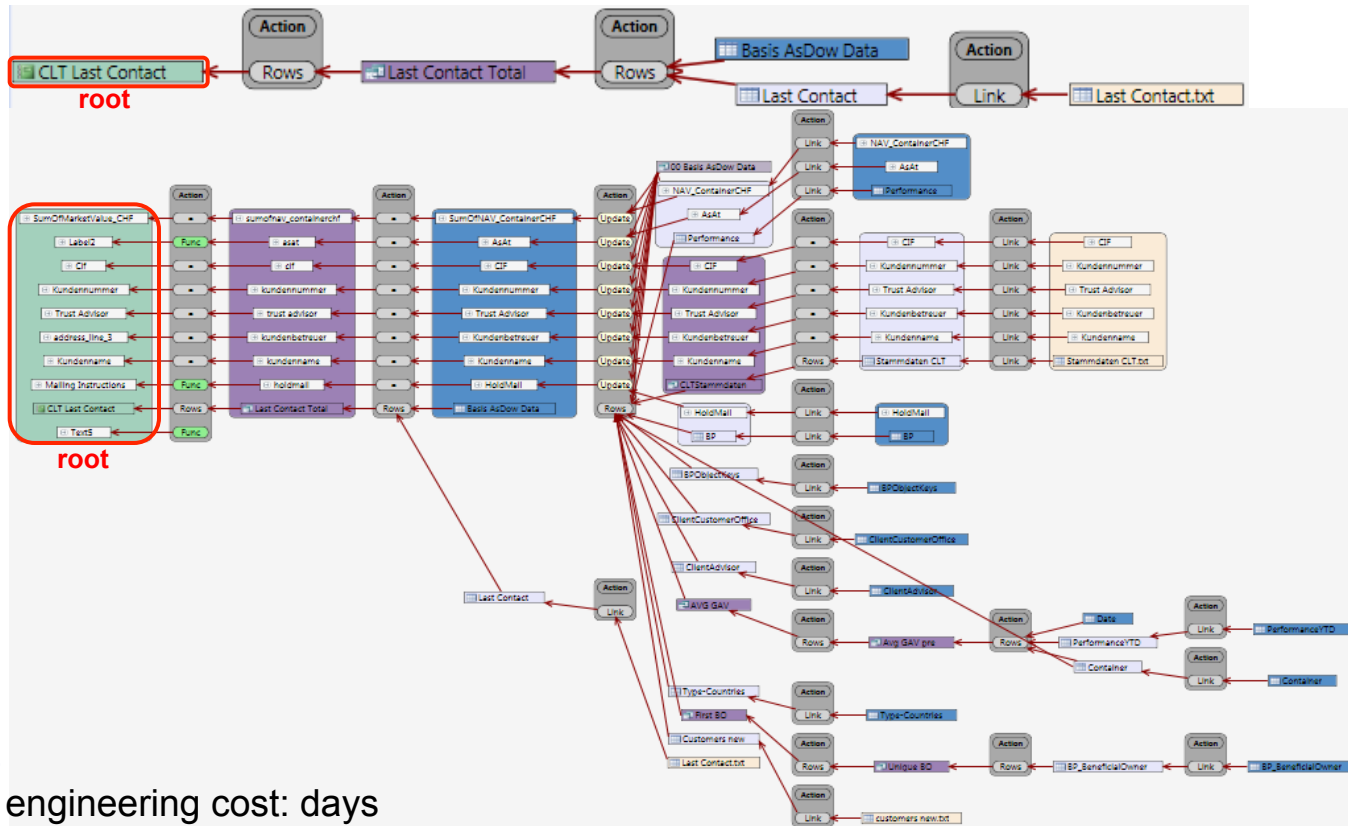
- extract all VB code from the Access reports
- perform a full **syntactic and control-flow analysis**
 - find VB code that writes to report fields
 - find how that code is called (call analysis)
 - trace back values of output field names to SQL field names using **VB symbolic execution**

easy
relatively easy
very complex!

Solution - putting it all together

Access Analyzer (SolidAA)

- end-to-end dataflow analysis across entire reporting platform
- answers question: “Where does this (report) data come from?”
- fully handles any MS Access / SQL database



Benefits

- reverse engineering cost: days
- learning cost: days
- solution cost: ~6 months
- client estimated savings: **~920 KEUR**

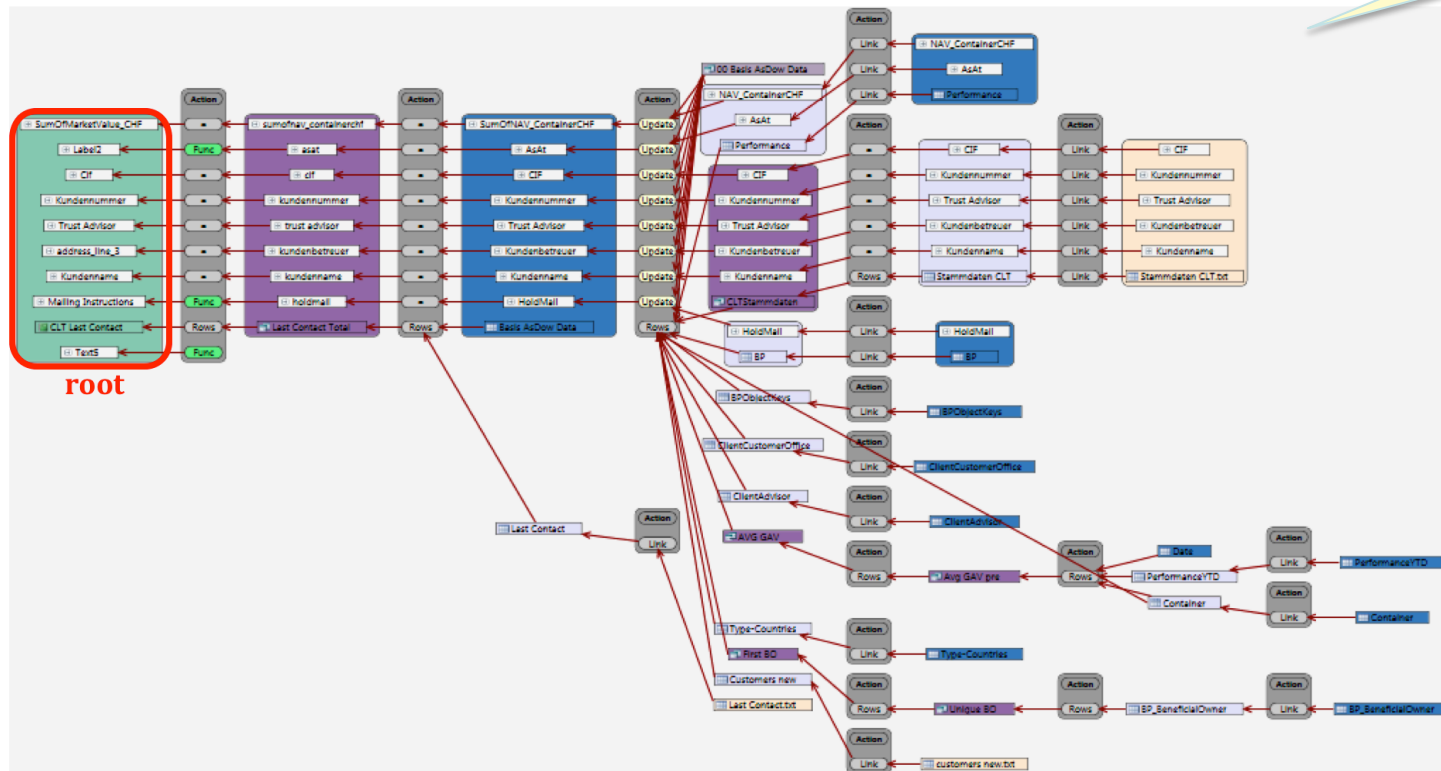
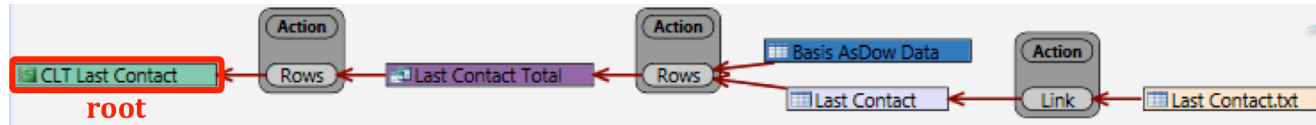
Solution refinements

Level of detail

- show dependencies at form and table level only (overview)
- expand to show dependencies at form-field and table-field level (detail)
- detect “parallel paths” of data processing (desirable w.r.t. architecture)

Overview

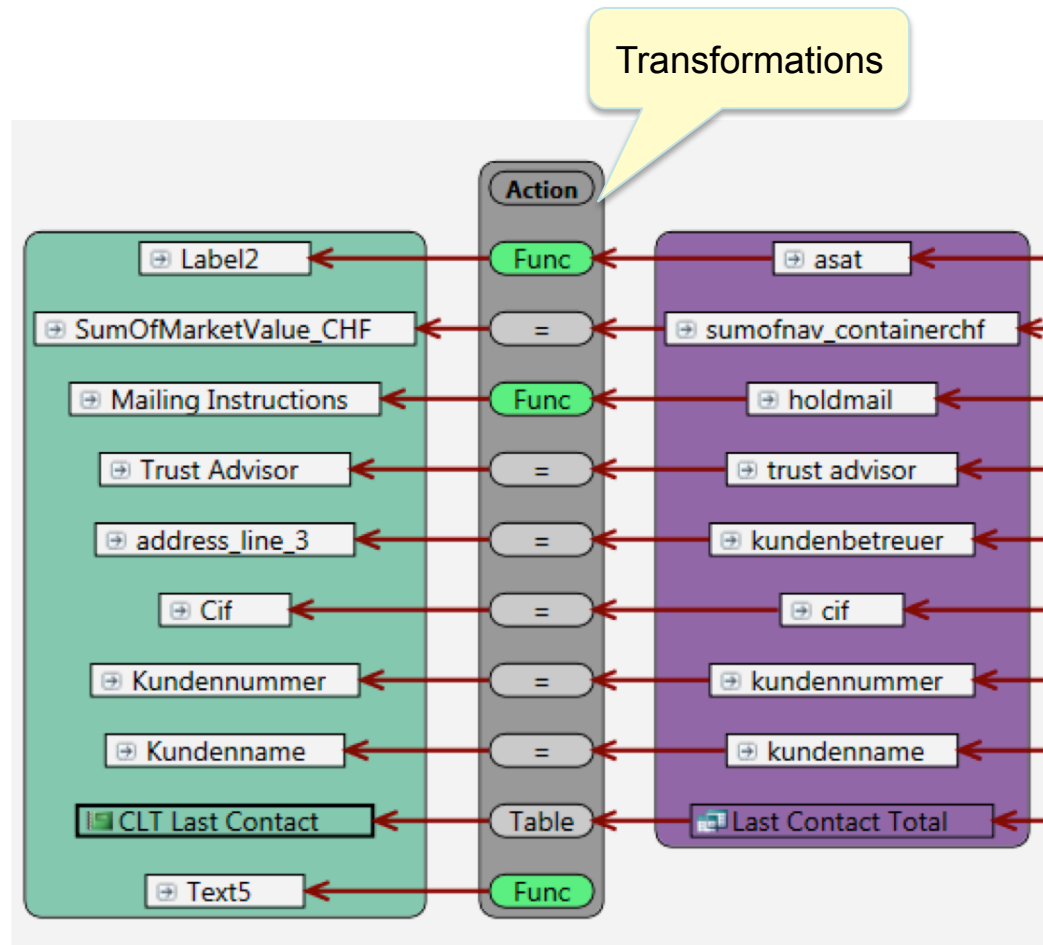
Detail



Solution refinements (cont.)

Level of detail (cont.)

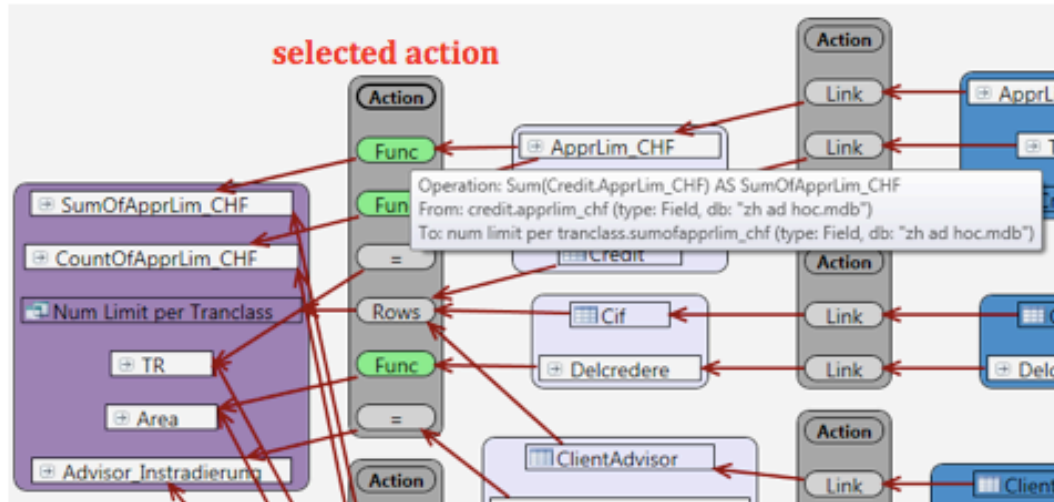
- show how the data is **transformed** from input to output
 - copy, SQL function call, references (e.g. table links), ...



Solution refinements (cont.)

Links to code (cont.)

- show which bits of code (SQL, VB, Access, ...) are responsible for each data link
- useful for fine-grained detail and understanding of the business logic



dependency view

action code

```
Code View
File: C:\Users\Public\SolidSource\serialized\Selection\ZH A Line: 3 Open
query: Num Limit per Tranclass
type: 0
connect:
sql: SELECT DISTINCT Left[[Delcredere],4] AS Area, ClientAdvisor.Advisor
FROM Credit INNER JOIN (Cif LEFT JOIN ClientAdvisor ON Cif.Pid=Client
GROUP BY Left[[Delcredere],4], ClientAdvisor.Advisor_Instradierung, Cred
HAVING (((Left[[Delcredere],4]) In ("XRBA","XRBD","XRBE"))))
ORDER BY Left[[Delcredere],4], ClientAdvisor.Advisor_Instradierung, Cred
```

code view

Application: Build Optimization

1. Context

- major embedded software company (NASDAQ 100)
- industrial 17.5 MLOC code base of C code
- modified daily by >500 developers worldwide

2. Problem

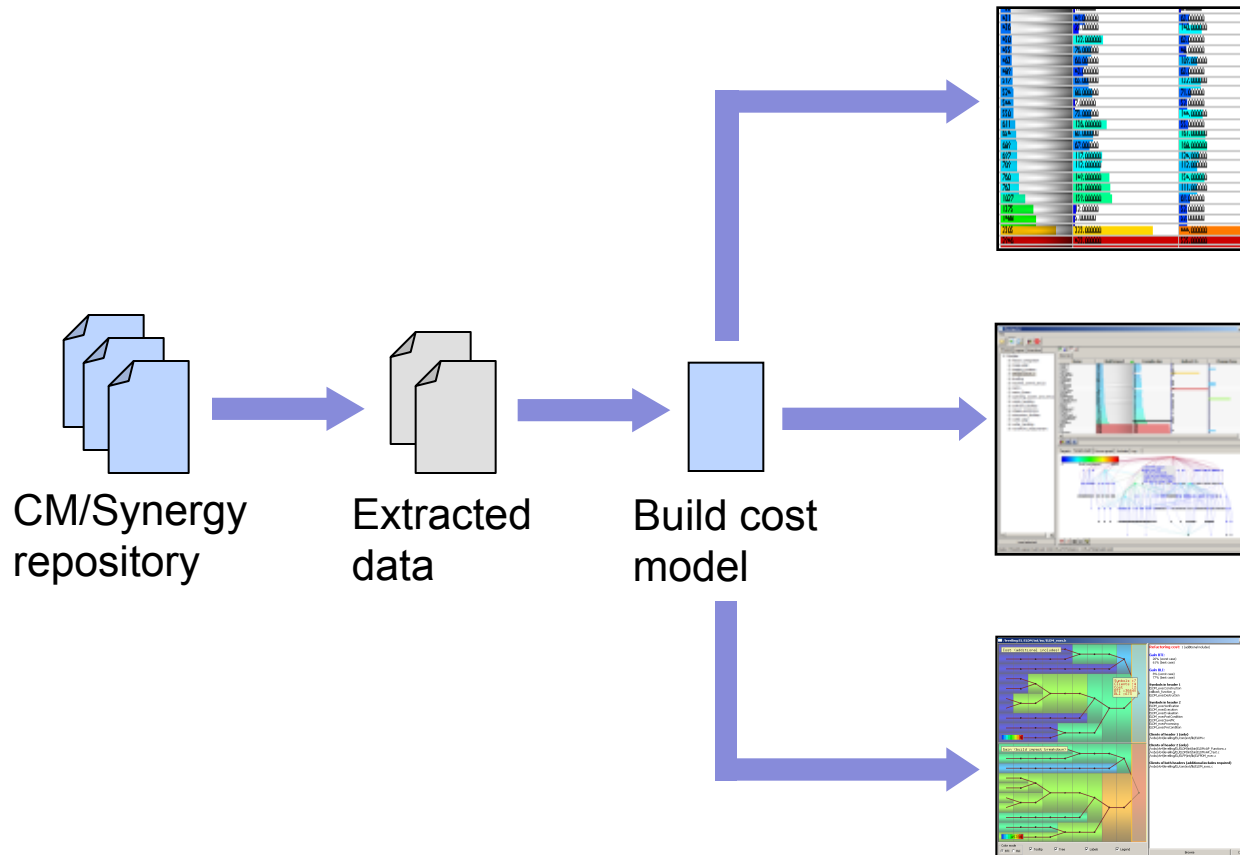
- high build time (>9 hours)
- modifying a header causes **very long recompilations**
- testing becomes very hard; perfective maintenance (refactoring) nearly **impossible**

3. Questions

- why is the build time so **long**?
- what **impact** has a code change on the build time?
- how is a change impact **spread** over the entire code base?
- how to **refactor** the code to improve modularity and build time?

Build Optimization

Three analyses – three tools in a unified toolset



TableVision tool

Build process analysis

- why is the build **slow**?

INavigator tool

Dependency analysis

- how does a code **change** affect build time?

IRefactor tool

Refactoring analysis

- how to **rewrite** code to improve build time?

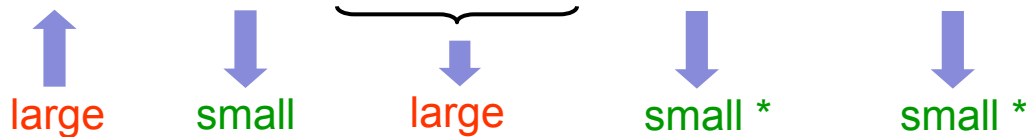
Question 1: Why is the build slow?

- measure build time using UNIX tools `time(x)`

build time	CPU time
792.500000	8.200000
767.599976	6.200000
745.900024	39.500000
722.299988	7.300000
719.700012	14.000000
694.500000	4.200000
688.299988	9.300000
673.299988	4.600000

large small

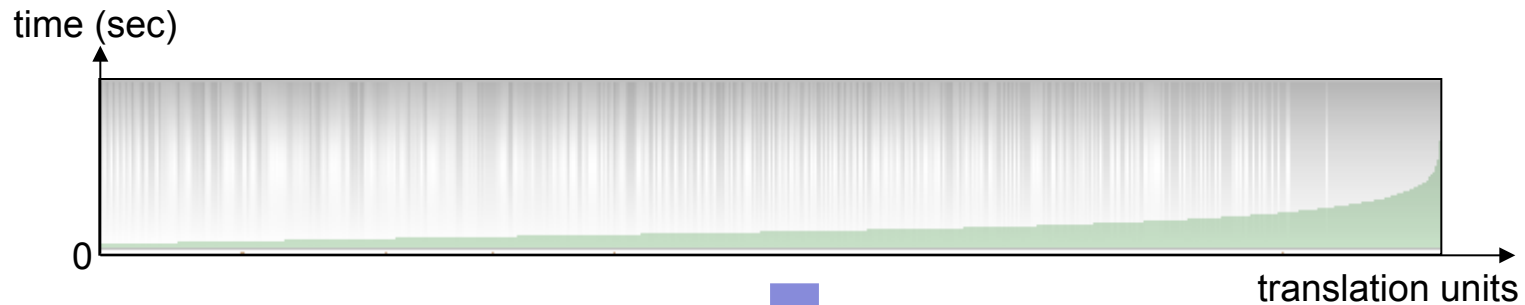
- build time = CPU + I/O + network + paging + other processes



* assume no other CPU-intensive processes besides compilation

Question 1: First Steps

- simple histogram of build time



Build time depends **significantly** on the translation unit!

A useful build cost model must consider the **per-unit build cost** and **not only** the number of translation units

Build Cost Model - First Attempt

Build cost

“how much it costs to build a file”

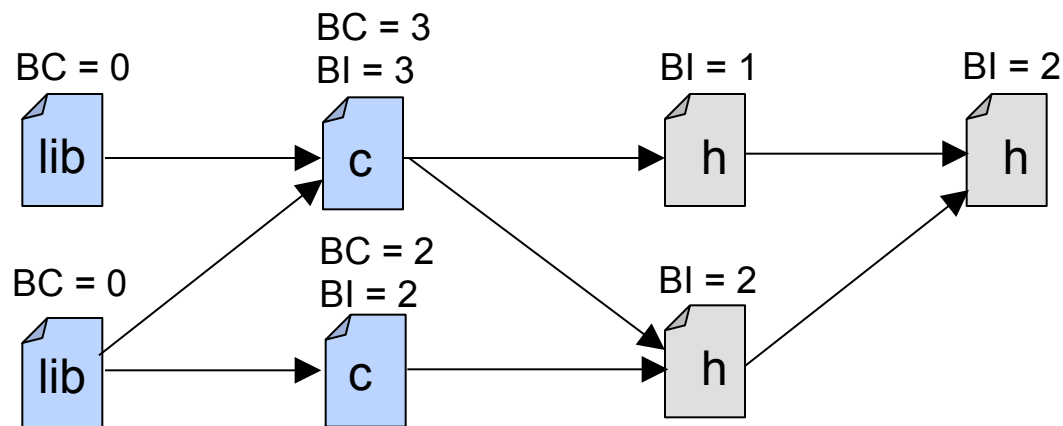
- sources: number of source LOC + (in)directly included headers
- binaries: negligible (linking is cheap)
- headers: zero (headers don't get compiled)

Build impact:

“how much it costs to rebuild the system when a file is modified”

- sources: build cost of the source itself
- headers: number of sources using that header

Example



* both application and system headers are considered

Build Cost Model - Validation

sorted on
model's impact

headers	model's impact	time build measurements
AM/com/ext/inc/AM.h	3302	304316
DNDM/com/int/inc/DNDMxADMIN.h	6396	593680
DD/com/ext/inc/DDXA.h	6399	516282
IM/com/ext/inc/IM.h	6432	595957
LU/com/ext/inc/LUGB.h	6699	630971
UL/com/ext/inc/ULXAerr.h	7084	780792
HW/com/ext/inc/HW.h	7376	645346
SM/com/ext/inc/SMXA.h	9449	976110
SW/com/ext/inc/SW.h	9948	998233
LU/com/ext/inc/LU.h	17864	1682591
MI/MITF/int/inc/MIFREV.h	18380	1273473
DD/com/ext/inc/DDXA_types.h	18380	1273474
ER/com/ext/inc/ERXAerror_types.h	18380	1273473
GW/com/ext/inc/GWXA.h	18380	1273473
MO/MPPC/int/inc/MO.h	18380	1273473
ZS/com/int/inc/ZSUM_umf_internal.h	18380	1273474
DM/com/int/tst/DMTP.h	18380	1273473
MO/com/ext/inc/MOERR.h	18380	1273473
CN/com/ext/inc/CNDD.h	18380	1273473
CN/com/ext/inc/CN.h	36760	2546947
MI/MIHP/ext/inc/PCHPIF.h	36760	2546947
PL/com/ext/inc/PLXA.h	36760	2546947
MG/MGTAB/ext/inc/MGTAB.h	36760	2546947
ER/ERLO/ext/inc/ERXA.h	36760	2546948
MI/com/ext/inc/MI.h	36760	2546947
MI/MIHP/ext/inc/PCHP.h	36760	2546947
OO/com/ext/inc/OOXA.h	55140	3820421
MG/com/ext/inc/MG.h	55140	3820421
CN/com/ext/inc/CNXA.h	55140	3820421
TH/com/ext/inc/THXA.h	55140	3820421

↑ low-impact headers

12th highest-impact header (reality)
classified as 21st (model)

↓ high-impact headers

- model is close to reality but not perfect
- deviations are important!

Build Cost Model - Refinement

Build cost

“how much it costs to build a file”

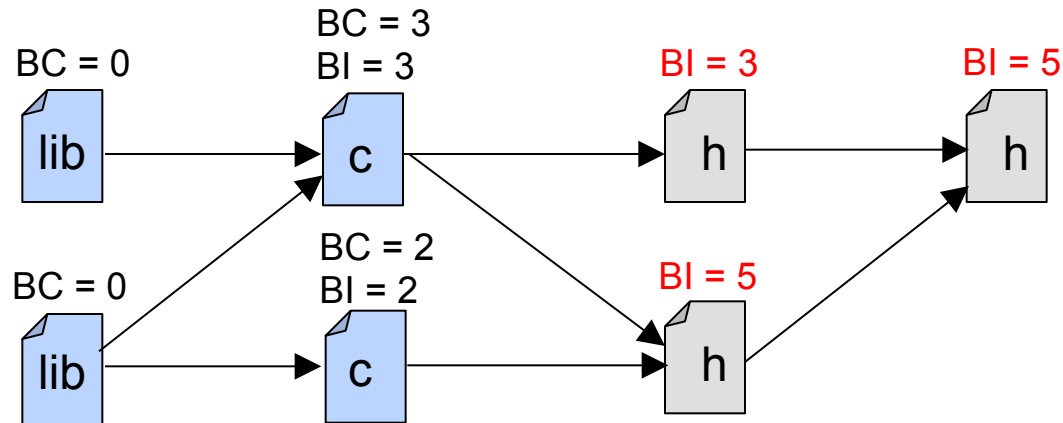
- sources: number of (in)directly included headers
- binaries: negligible (linking is cheap)
- headers: zero (headers don't get compiled)

Build impact:

“how much it costs to rebuild the system when a file is modified”

- sources: the build cost of the source itself
- headers: **sum of build costs of all sources** including header (in)directly

Example



Build Cost Model 2 - Validation

sorted on
first model's impact

headers	refined model's impact	first model's impact	time build measurements
AM/com/ext/inc/AM.h	1141036	3302	304316
DNDM/com/int/inc/DNDMxADMIN.h	2231400	6396	593680
DD/com/ext/inc/DDXA.h	1717353	6399	516282
IM/com/ext/inc/IM.h	2240332	6432	595957
LU/com/ext/inc/LUGB.h	2217519	6699	630971
UL/com/ext/inc/ULXAerr.h	2067954	7084	780792
HW/com/ext/inc/HW.h	2471980	7376	645346
SM/com/ext/inc/SMXA.h	2603546	9449	976110
SW/com/ext/inc/SW.h	3433500	9948	908233
LU/com/ext/inc/LU.h	5913384	17864	1682591
MI/MITE/int/inc/MIFREV.h	4841045	18380	1273473
DD/com/ext/inc/DDXA_types.h	4841045	18380	1273474
ER/com/ext/inc/ERXAerror_types.h	4841045	18380	1273473
GW/com/ext/inc/GWXA.h	4841045	18380	1273473
MO/MPPC/int/inc/MO.h	4841045	18380	1273473
ZS/com/int/inc/ZSUM_umf_internal.h	4841045	18380	1273474
DM/com/int/tst/DMTP.h	4841045	18380	1273473
MO/com/ext/inc/MOERR.h	4841045	18380	1273473
CN/com/ext/inc/CNDD.h	4841045	18380	1273473
CN/com/ext/inc/CN.h	9682090	36760	2546947
MI/MIHP/ext/inc/PCHPIF.h	9682090	36760	2546947
PL/com/ext/inc/PLXA.h	9682090	36760	2546947
MG/MGTAB/ext/inc/MGTAB.h	9682090	36760	2546947
ER/ERLO/ext/inc/ERXA.h	9682090	36760	2546948
MI/com/ext/inc/MI.h	9682090	36760	2546947
MI/MIHP/ext/inc/PCHP.h	9682090	36760	2546947
OO/com/ext/inc/OOXA.h	14523135	55140	3820421
MG/com/ext/inc/MG.h	14523135	55140	3820421
CN/com/ext/inc/CNXA.h	14523135	55140	3820421
TH/com/ext/inc/THXA.h	14523135	55140	3820421



low-impact headers

refined model classifies outlier correctly

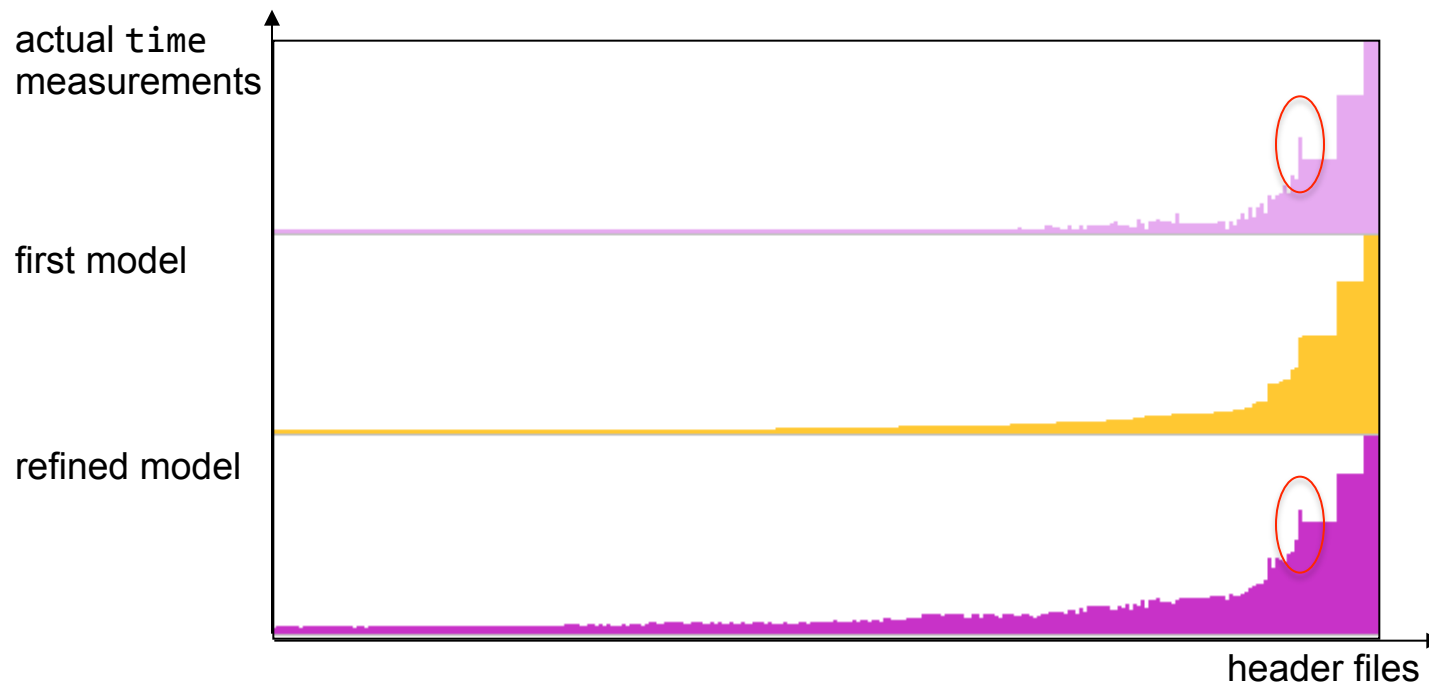


high-impact headers

- refined model delivers same header-order (in terms of impact) as actual measurements

Build Cost Model 2 - Validation

Let's look at the whole picture

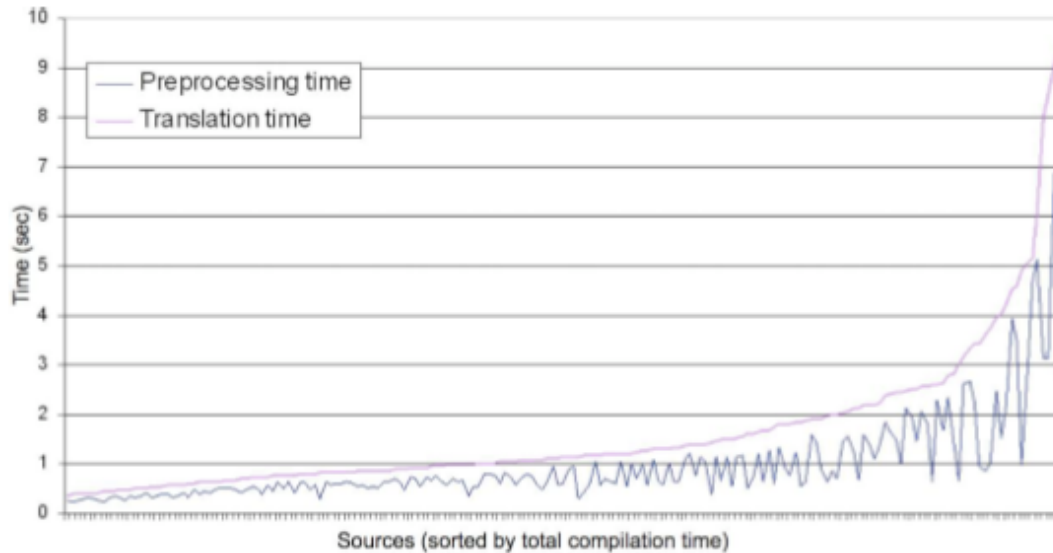


- refined model nicely matches reality, including subtle 'outliers'
- why is this so? (see next slide)

Build Cost Model 2 - Validation




Analyze deeper:

- compilation cost dominated by I/O (preprocessing headers)
- I/O cost dominated by file opening/closing on this platform
- hence the justification of impact = # totally opened headers



Conclusions

To reduce build time, we should:

- either massively accelerate network → highly **costly** / **complex** 
- reduce per-header build impact → header impact analysis 
- reduce impact of change on build time → header refactoring 

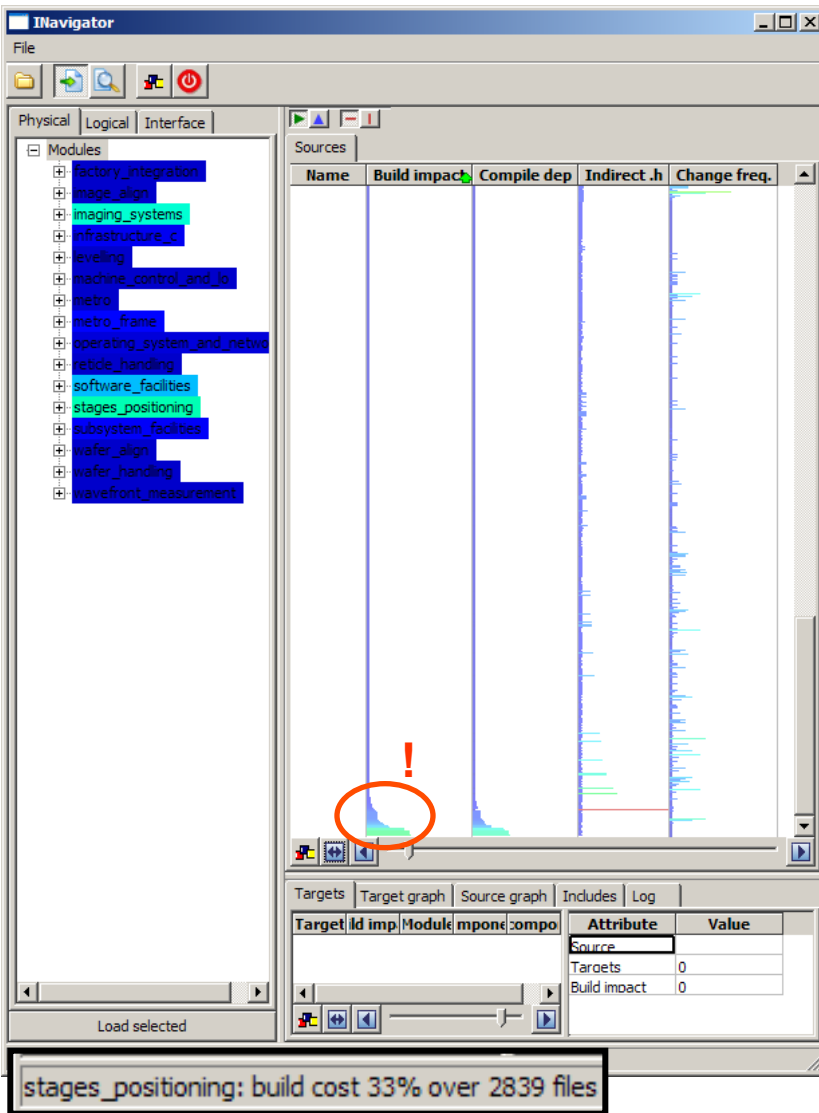
Question 2: System-wide impact analysis

1. Find subsystems are expensive to build
2. For a subsystem, find headers have high build impact
3. Zoom in to highest impact headers
- 4,5. For a high-impact header, see how its impact spreads over sources
6. For a header, see its cost breakdown over its include-set

The screenshots illustrate the following steps:

- Step 1:** The Navigator interface shows a list of modules on the left and a table of sources with columns for Name, Build impact, Compile dep, Indirect.h, and Change freq.
- Step 2:** A zoomed-in view of the source table, highlighting a specific header with a high build impact.
- Step 3:** A detailed view of the selected header's properties, including its name, build impact, and compile dependencies.
- Step 4:** A target graph showing the impact of the selected header on other sources in the system.
- Step 5:** A source graph showing the build cost breakdown of the selected header over its include-set.
- Step 6:** A complex network graph showing the build cost breakdown of the selected header over its include-set, with nodes representing sources and edges representing dependencies.

Subsystem-level impact analysis



Method

- color system tree by cost (blue=low, red=high)
- select desired subsystem
- right panel shows **build impact** for each header / source in that subsystem

Findings

- most headers have a low build impact
- however, a few have a very high impact
- touching those incurs a **high build cost!**

- ➔ because they are used in many sources
- ➔ because they include many headers

Question 3: How to reduce the build cost?

- OK, we have a high-impact header h : how easy it to **reduce** that impact?
- visualize the build cost distribution of h over the sources which use it

Case 1: easy refactoring

- build cost spread **unevenly** over the targets including selected header h
- to decrease cost due to h , we only need to change a few **targets**

WSPFBTyp.h	571406	1763	222	0
WSSUPTyp.h	586690	1830	221	0
STIOTyp.h	618630	1962	220	0
WSXATyp.h	657773	2094	221	0
WXXATyp.h	670105	2177	219	0
WSXAxCHUCKTyp.h	761159	2488	1	0
PCHP.h	2318254	9058	1	3
MIXATypi.h	2318254	9058	218	0
ML.h	2318254	9058	218	3
MIXA.h	2318254	9058	218	0
PCHPIF.h	2318254	9058	218	3
MIXATyp.h	2318254	9058	218	0
OMTT_satoolio.h	2318254	9058	218	0
MIFmw.h	2318254	9058	1	0

Target	Build imp	Module	Component	Component
THMain CC.c	307	software facilities	TH	THMA
TSDM.c	419	stages positioning	TS	TSDM
WS_error.c	409	stages positioning	WS	com
TSPEXAP_sp.c	398	stages positioning	TS	TSPE
TSDAAP_analyse.c	394	stages positioning	TS	TSDA
TDXTxWH.c	391	stages positioning	TD	com
TSDAAP_main.c	388	stages positioning	TS	TSDA
TSSM_exec.c	386	stages positioning	TS	com
SWXBxAP_Measure.c	383	stages positioning	SW	SWXB
TSREXAP_tt_callb.c	383	stages positioning	TS	TSRE
TSDAAP_exec.c	383	stages positioning	TS	TSDA
SWLFXAP_sp.c	382	stages positioning	SW	SWLF

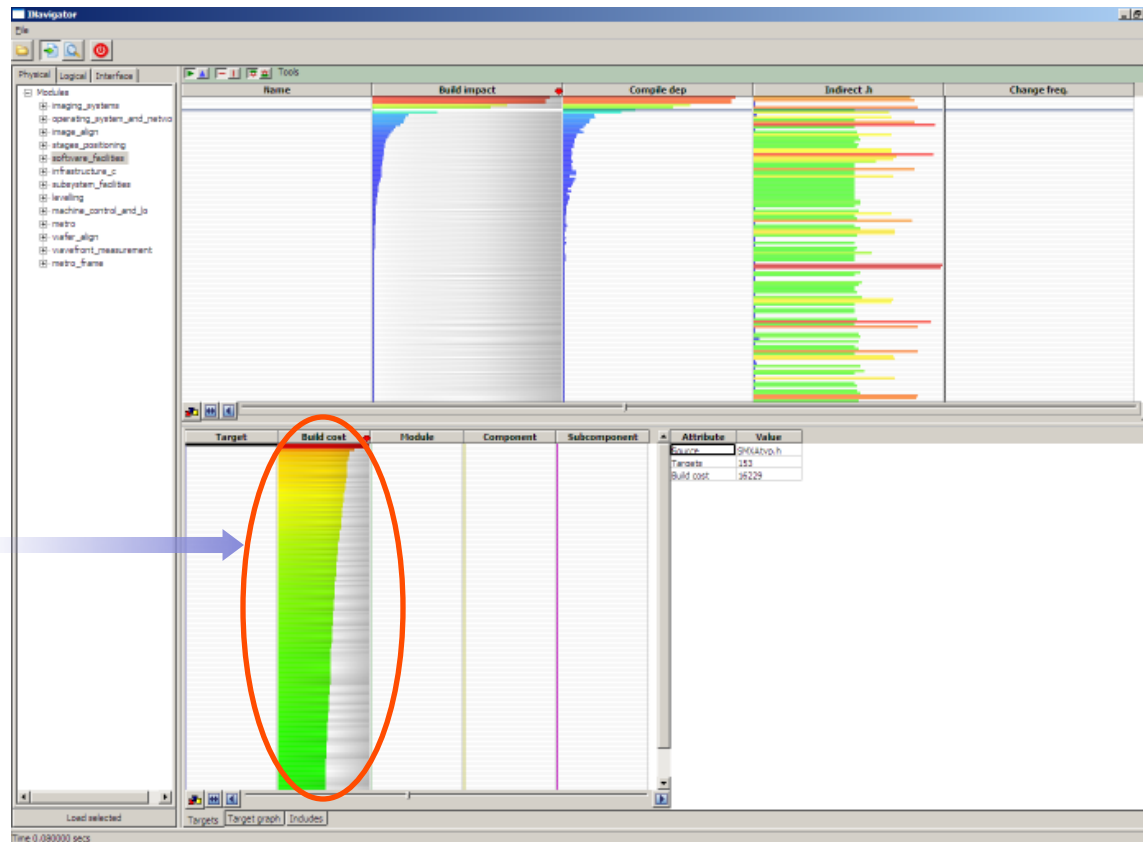
build impact of h is located mainly in **one** single place!

Target	Build imp	Module	Component	Component
THMain CC.c	307	software facilities	TH	THMA
TSDM.c	419	stages positioning	TS	TSDM
WS_error.c	409	stages positioning	WS	com
TSPEXAP_sp.c	398	stages positioning	TS	TSPE
TSDAAP_analyse.c	394	stages positioning	TS	TSDA
TDXTxWH.c	391	stages positioning	TD	com
TSDAAP_main.c	388	stages positioning	TS	TSDA
TSSM_exec.c	386	stages positioning	TS	com
SWXBxAP_Measure.c	383	stages positioning	SW	SWXB
TSREXAP_tt_callb.c	383	stages positioning	TS	TSRE
TSDAAP_exec.c	383	stages positioning	TS	TSDA
SWLFXAP_sp.c	382	stages positioning	SW	SWLF

Refactoring analysis

Case 1: difficult refactoring

- build cost spread **evenly** over the targets including selected header h
- to decrease cost due to h , we need to change almost **all targets**



← selected high-impact header

build impact of h →

Refactoring analysis - Refinement

- not all headers change equally often (e.g. system headers)
- new metrics:
 - build impact * change frequency
 - impact distribution: impact (%) of a header contained in the 10% most expensive of its targets
- easy & quick to use

flat, low distribution (~15%): impact is spread uniformly over all targets. Hence, we cannot improve by refactoring a few targets

skewed distribution (50%): half of impact is concentrated in 10% most expensive targets. Hence, refactoring these is an interesting option

Name	Impact/targets	Cost	Change freq.	Impact*freq.	Targets*freq.	Distribution
SN.h	10368 33	266	2	20736	66	13
MGH_CB_HSI.h	21592 67	1	1	21592	67	10
MSO_OB1_Pt.h	24283 77	270	1	24283	77	11
MSO_OB3.h	24283 77	269	1	24283	77	11
MSEV.h	24283 77	266	1	24283	77	11
JAHW.h	12440 45	265	2	24880	90	11
MSFERER.h	29984 97	270	1	29984	97	11
MSFERER_errorc	30255 98	3	1	30255	98	11
MSO_SI_LIN.h	17095 51	279	2	34190	102	11
MSO_AI_LIN_AO	17095 51	280	2	34190	102	11
MSO_AI_LIN.h	17095 51	279	2	34190	102	11
MSO_PERF_SETT	17370 52	272	2	34740	104	11
MSO_CHK_KIN_E	17370 52	272	2	34740	104	104
MSO_PERF_SERV	17370 52	272	2	34740	104	104
MSO_CHK_FORC	17373 52	273	2	34746	104	11
MSO_CHK_COLL	17373 52	273	2	34746	104	11
MSO_CHK_STAB	17373 52	273	2	34746	104	11
MGH_Pt.h	19043 66	272	2	38086	132	11
MSACCN.h	40638 142	278	1	40638	142	11
MSFERER.h	42592 149	270	1	42592	149	10
MSET_errorbases	39705 132	2	2	79410	264	11
MSO.h	40638 142	278	2	81276	284	11
MSOsubact.h	51371 180	266	3	154113	540	10
MGVd.h	44894 157	269	5	224470	785	10
MGPpCdef.h	84259 288	4	3	252777	864	11
MSxERRORBASE	82075 316	1	4	368300	1264	11
MSL1.h	289538 860	265	3	869514	2580	12
DMTP.h	82686 17459	264	1	82686	17459	

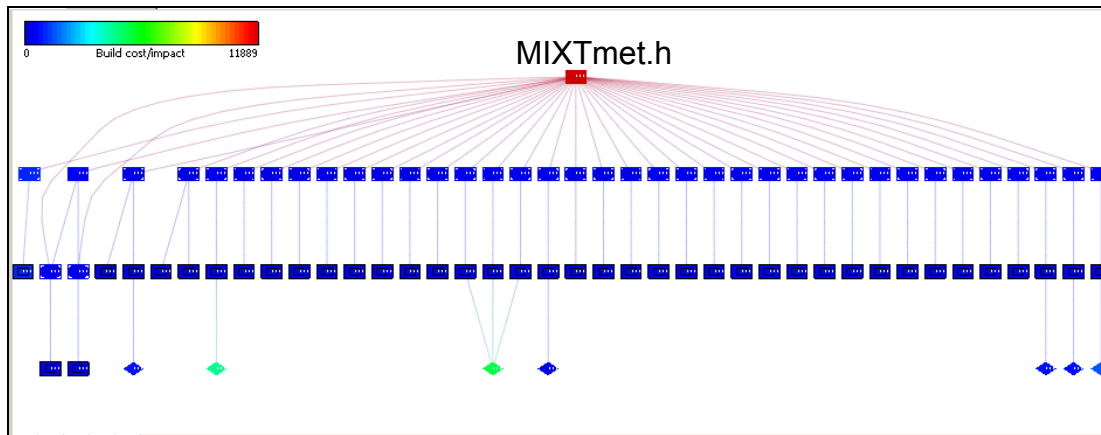
sorted by impact*change

Name	Impact	Targets	Cost	Change freq.	Impact*freq.	Targets*freq.	Distribution
HWNAP_exe	299	4	259	4	1196	4	0
WHFAT_c	299	4	259	4	1196	4	0
WHFAPS_sbo.c	300	4	300	4	1200	4	0
HWNAP_exe.c	300	4	300	4	1200	4	0
HWNAP_exe.c	300	4	300	4	1200	4	0
WH_RQ_QUEUE.c	301	4	301	4	1204	4	0
WHDOCs_dpr.h	1208	4	264	1	1208	4	25
WHDOCs_imp.c	1212	4	271	1	1212	4	25
WHOT_test.c	306	4	306	4	1224	4	0
WHAPPA_Jul.h	812	2	287	2	1624	4	51
WHESU_dpr.c	307	4	307	4	1228	4	0
WHAPLM_clear.c	308	4	308	4	1232	4	0
WHAPLM_irq.c	308	4	308	4	1232	4	0
WHAPLM_irq.c	308	4	308	4	1232	4	0
WHAPLM_irq.c	308	4	308	4	1232	4	0
WHADAP_OV_ovl232	272	4	272	4	1088	4	27
WHAPRO_test.c	309	4	309	4	1236	4	0
WHAPPA_unhook16	2	2	265	2	1232	4	51
WHAPPA_irq.c	1125	4	308	4	1235	4	25
WHAPLM_event1239	4	4	275	1	1239	4	0
WHDOCs_c	310	4	310	4	1240	4	0
WHAPLM_irq.c	312	4	312	4	1248	4	0
WH_RQ_event.c	313	4	313	4	1252	4	0
WHAPRO_move	315	4	315	4	1260	4	0
WHAPRO_irq.c	317	4	317	4	1268	4	0
WHADAP_LW_Jo	1272	4	270	4	269	4	0
WHAPRO_remove1314	4	4	278	1	1272	4	0
WHAPRO_irq.c	317	4	317	4	1268	4	0
WHAPRO_retr	1304	4	278	1	273	4	0
WHAPRO_irq.c	317	4	317	4	1268	4	0
WHAPPA_irq.c	1417	5	269	5	269	4	0
HWNAP.h	1420	5	275	5	275	4	0
HWNAP.h	1420	5	275	5	275	4	0
HWNAP.h	1425	5	275	5	275	4	0

Refactoring support

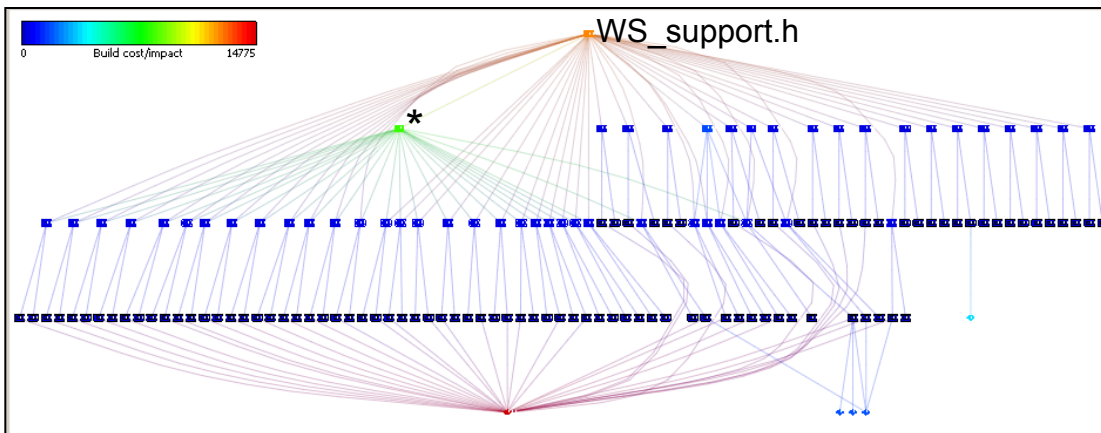
- OK, we found a high-impact header; how to decide a refactoring plan?
- show dependencies header → clients using hierarchical DAG layout

Example 1: MIXTmet.h, used by 38 sources, high impact



- build impact due to **direct header inclusion**
- hard to decrease via refactoring

Example 2: WS_support.h, used by 48 sources, high impact



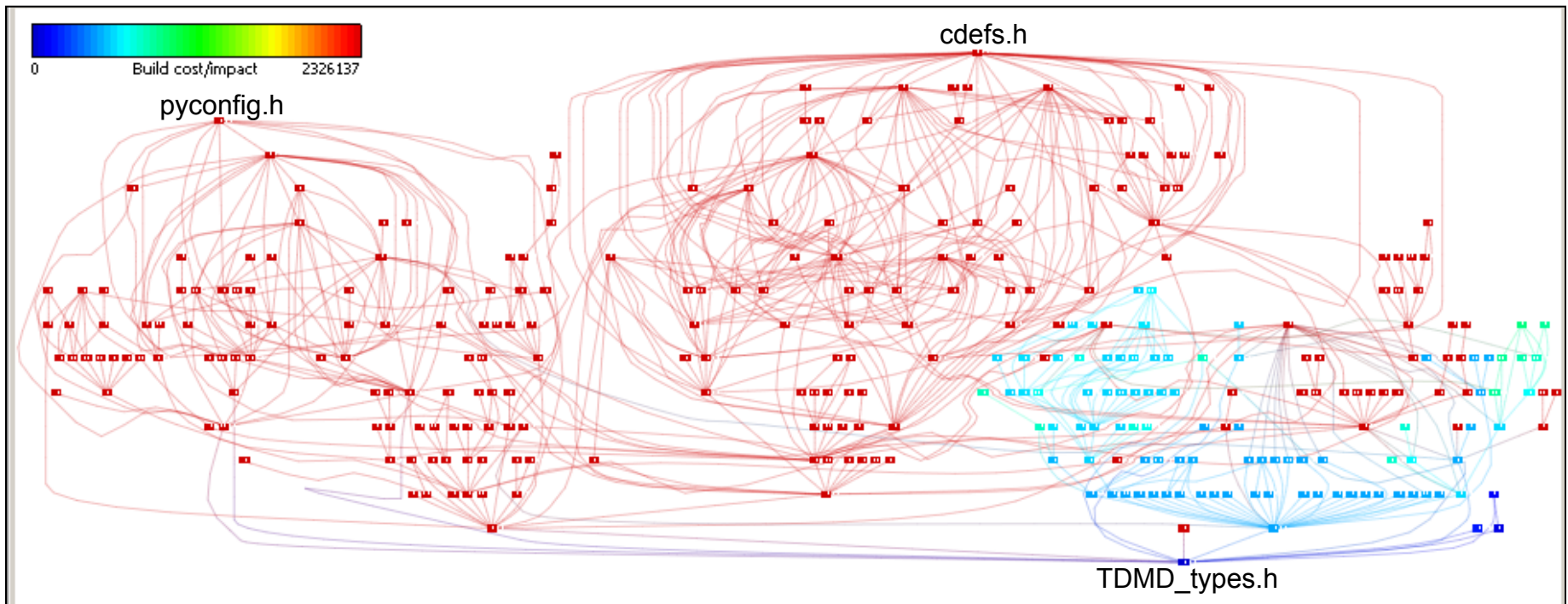
- build impact channeled via one **intermediate** header: WS_sim1_support.h (*)
- simpler refactoring may be possible

Refactoring support (2)

- say we want to include a header: is this potentially expensive?
- show header's own **include graph** colored by build impact

Example: TDMD_types.h, used by 30 sources

- not a high-impact header itself
- but it includes high-impact headers!
- hence using this header introduces potentially expensive changes



Refactoring support (3)

How much costlier becomes the system build if we add an #include?

- select a “source” header – the one in which we want to #include
- select a “destination” header – the one to be #included
- show the build cost increase

Example: What if we #include DNCHUI_chset.h in TDMD_types.h?

Name	Build impact	Compile dep	Indirect .h	Change freq.
IMAException.n	8483	26	3	1
WSSUP_NVtyp.h	8643	27	1	0
WSIOCO.h	8770	27	232	2
TSXTxPERFtyp.h	8866	25	1	0
MIXAx5Cmet.h	8912	28	219	0
WSAPIM_ER.h	8927	27	223	0
libWSPESS.so	9255	0	9255	0
SWSMAP.h	9354	27	245	2
libWSactuator.so	9357	0	9357	0
SWSMexception.h	9462	27	3	3
libTDXTxMD.so	9633	0	9633	0
TDMD_types.h	9633	30	312	1
SWSMxTCTyp.h	9713	28	220	0
SWSM.h	9713	28	55	0

source

Name	Build impact	Compile dep	Indirect .h	Change freq.
DNCHUI_chset.c	232	1	232	1
DNCHUI_chset.h	492	2	223	1
DNCHUItyp.h	721	3	220	0
DNCHexception.h	2560	11	2	0
DNCHxAPmet.h	716	3	221	0
DNCHxAPTtyp.h	1865	8	219	0
DNCHxDPTtyp.h	1853	8	110	0
DNCHUI_main.c	229	1	229	1
DNCHUI_mnwin.h	489	2	219	1

Current	Client	Include	Build impact	Indirect .h
YES	TDMD_types.h	-	9633	312
	TDMD_types.h	DNCHUI_chset.h	9783	317

target

build impact increases from 9633 to 9783, i.e. **1.5%**

Refactoring support (4)

- previous methods OK for manual header-by-header refactoring only

How to refactor a large system?

- system $S = \{f_i\}_i$, $S = Headers \cup Sources$
- header $h_i \in Headers = \{s_j\}_j$, $s_j \in Symbols$ (function declarations, variables, types, macros, ...)
- include relations
 - $inc : S \rightarrow \mathcal{P}(Headers)$, $inc(f) = \{h_i\} \Leftrightarrow f$ includes h_i
- symbol use relations
 - $use : Symbols \rightarrow \mathcal{P}(Headers)$, $use(s) = \{h_i\} \Leftrightarrow s$ is used by h_i
 - in typical systems, not all symbols $s_j \in h$ in a header are used **together**

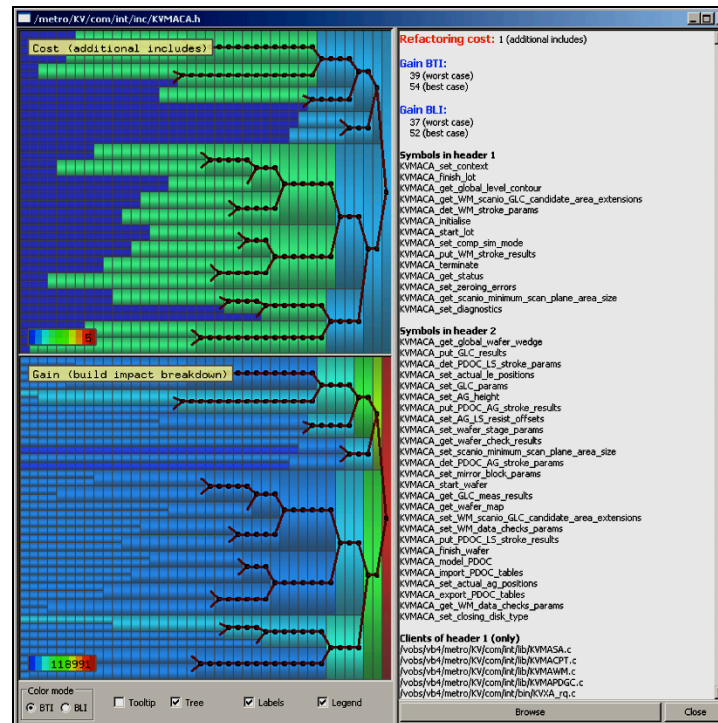
Automatic refactoring idea

- find high-impact header h (see last slides)
- split h into h_1, h_2 ; $h_1 \cup h_2 = h$ by putting symbols used together in same h_i
- recursively split h_1, h_2
- replace $inc(h)$ by $inc(h_1)$ and/or $inc(h_2)$

Refactoring support (4)

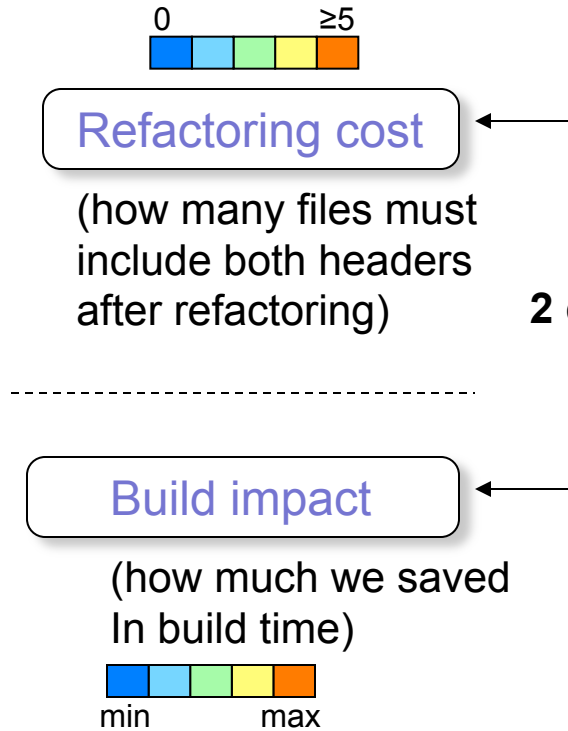
- intuitively: put symbols **often used together** in same header
- include newly created headers instead of original ‘monolithic’ ones
- why is this good
 - decrease build costs (by decreasing the included code size)
 - decrease build impact (by decreasing the number of included headers)

The IRefactor analysis tool

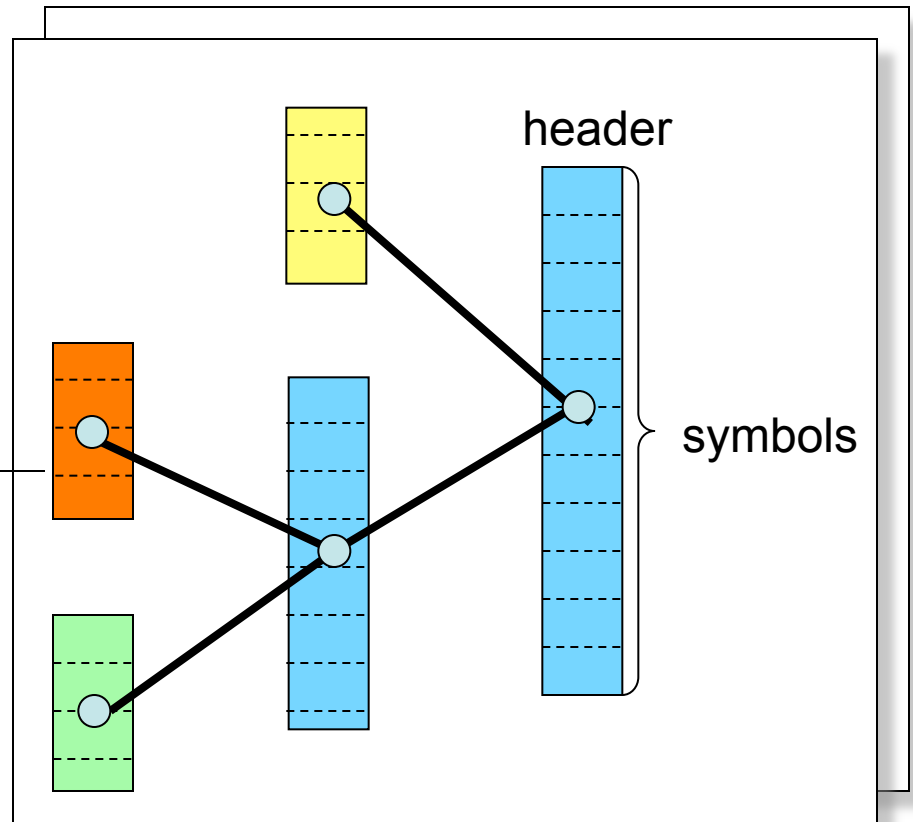


- suggests refactoring possibilities and shows **gained build impact**

Refactoring visualization



2 colors



Best refactoring candidates:

- low refactoring cost
- high build impact parents
- low build impact children

Refactoring cost



Build impact

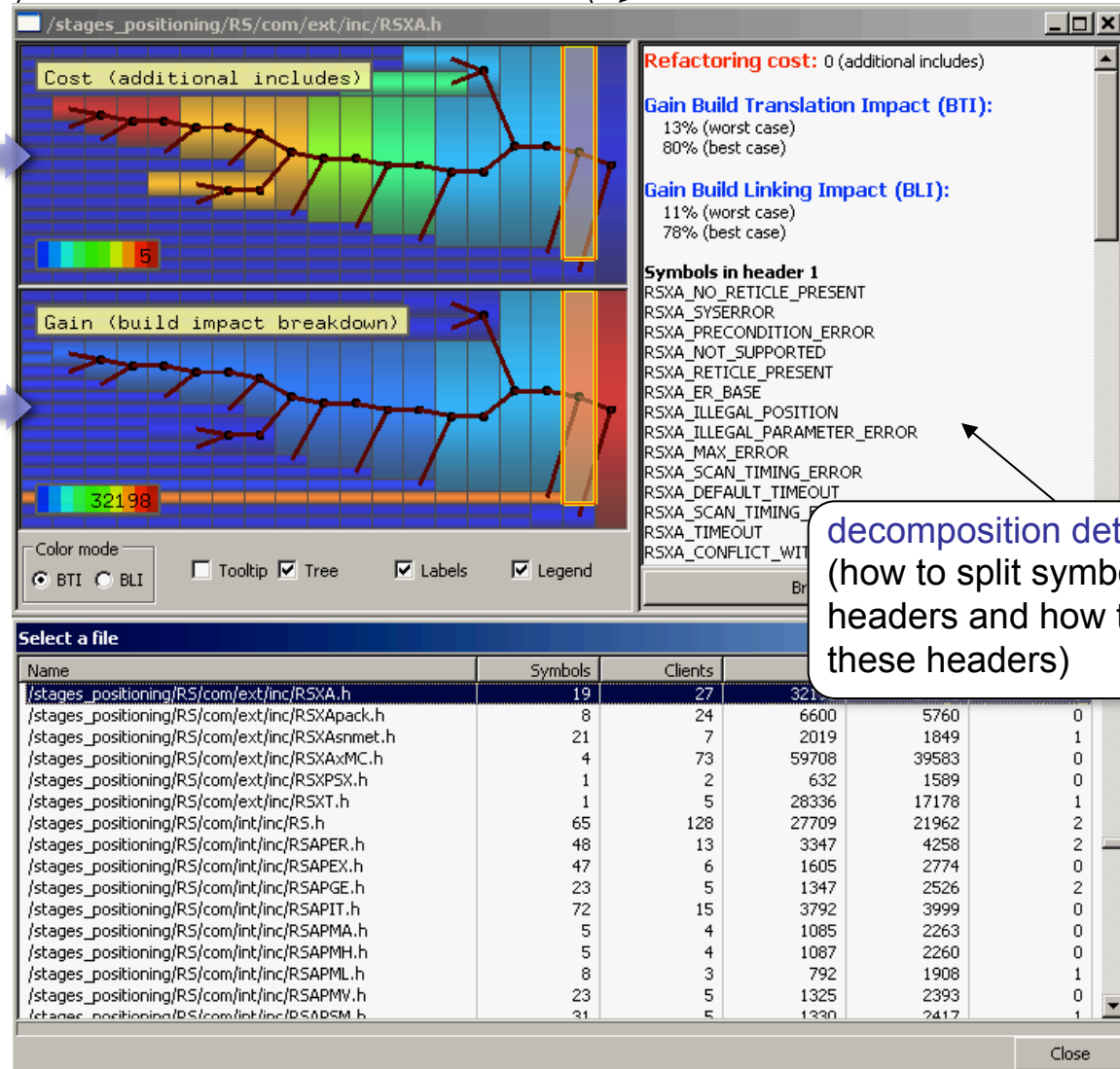
Refactoring visualization

suggested decomposition levels

header under analysis

Color: **refactoring cost**
(how many additional headers)

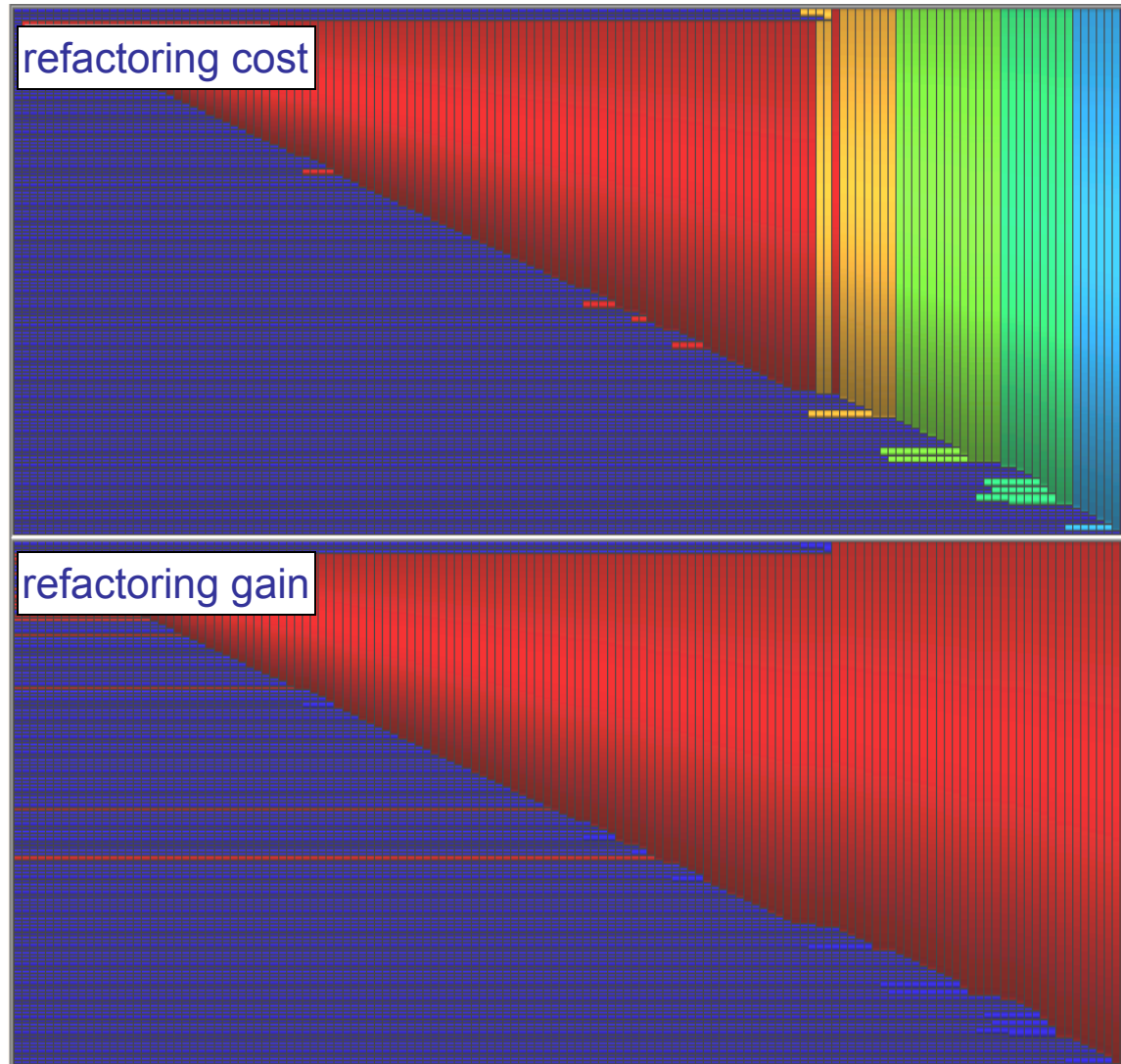
Color: **refactoring benefit**
(% build impact reduction)



Refactoring visualization

Example of bad candidate for header refactoring

As we gain benefits, we also increase costs

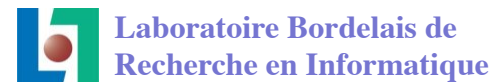


Collaborations & Customers

The Industry...



Research/Academia...

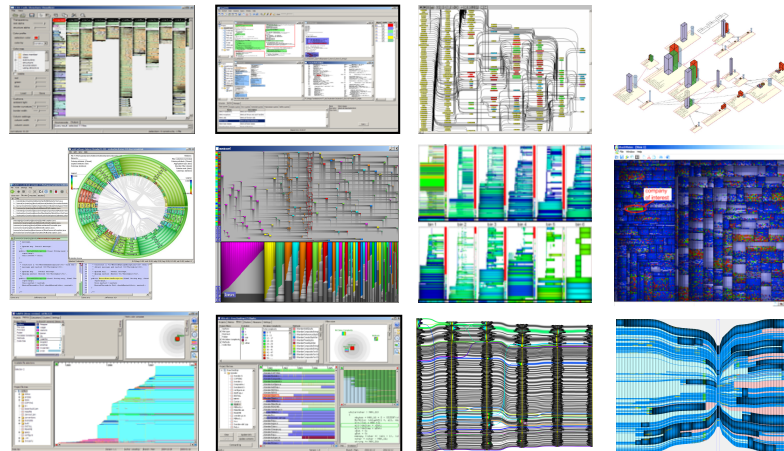


...Open Source



Conclusions - Software Visual Analytics

- **Provide insight in multidimensional correlations**
 - Program architecture, dependencies, metrics, development/testing effort, requirements, documentation, databases
 - Evolution of all these aspects in time
- **Added value**
 - Make the entire chain requirements..design..code visible and accountable
 - Assess software quality
 - Pinpoint hot-spots (where to invest effort)
 - Make sense of all that 'big data'



Thank you for your interest!

Alex Telea
a.c.telea@rug.nl