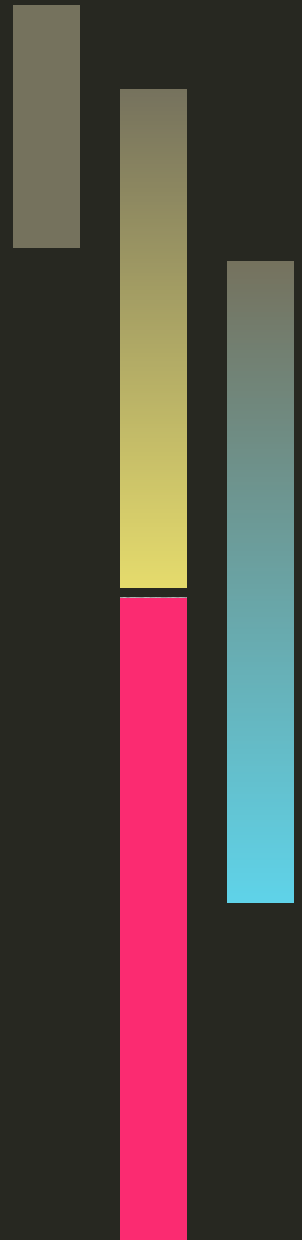


Reverse Engineering

Dr. Vadim Zaytsev aka @grammarware
UvA, MSc SE, 9 November 2015

Roadmap

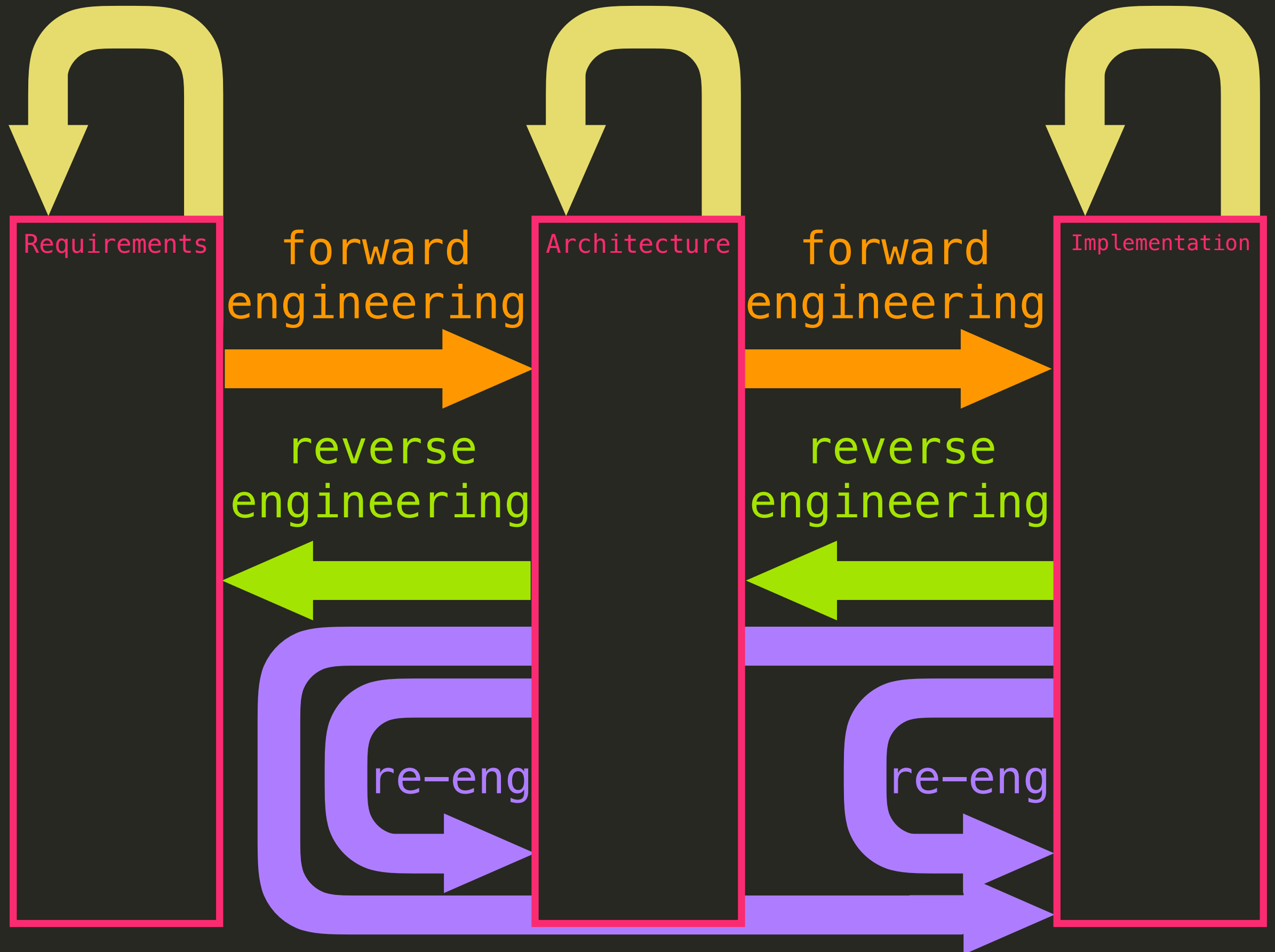
W44	Introduction	V. Zaytsev
W45	Metaprogramming	J. Vinju
W46	Reverse Engineering	V. Zaytsev
W47	Software Analytics	M. Bruntink
W48	Clone Management	M. Bruntink
W49	Source Code Manipulation	V. Zaytsev
W50	Legacy and Renovation	TBA
W51	Conclusion	V. Zaytsev



restructuring

restructuring

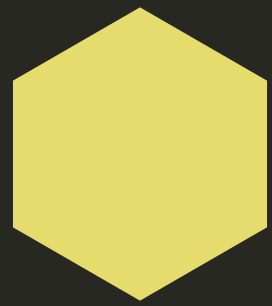
restructuring



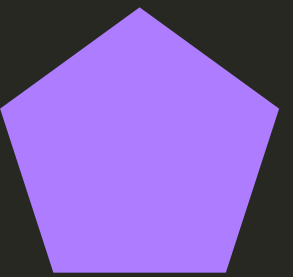
Objectives of reverse engineering

- * Cope with complexity
- * Generate alternate views
- * Recover lost info
- * Detect side effects
- * Synthesise higher abstractions
- * Facilitate reuse

Code Reverse Engineering

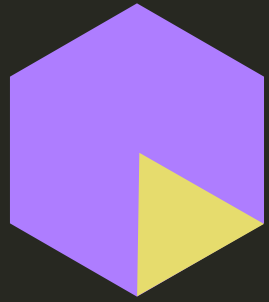


Code

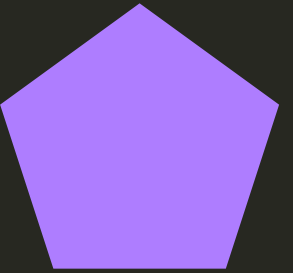


reverse engineering

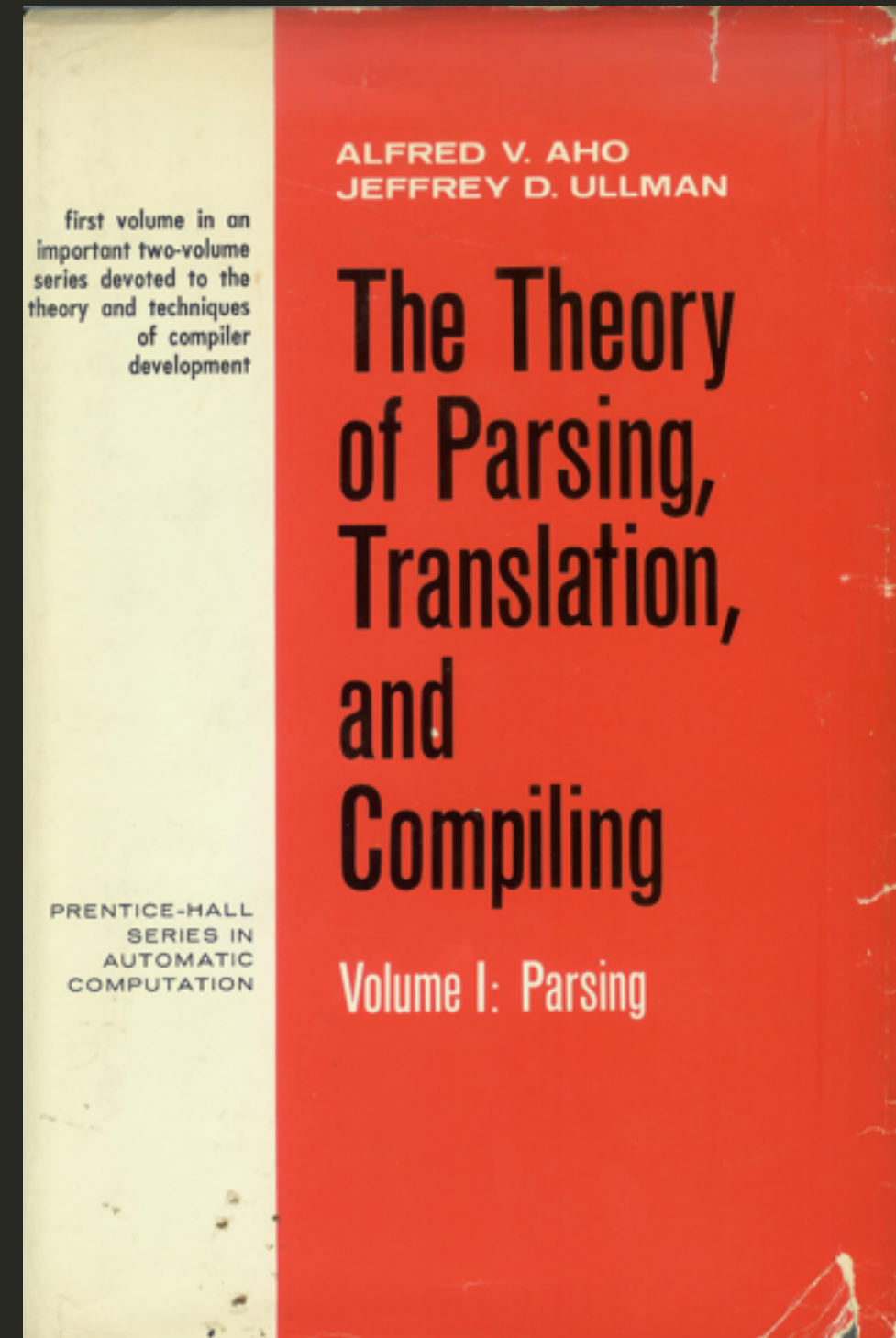
- * Parsing
- * Fact extraction
- * Slicing
- * Pattern matching
- * Decomposition
- * Exploration



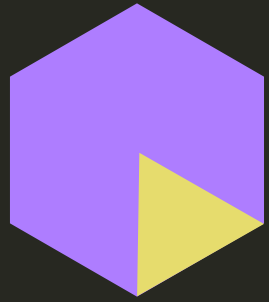
Parsing



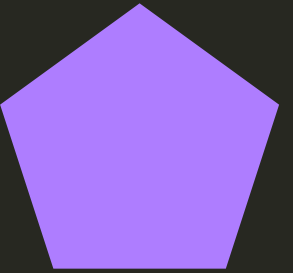
- * Well-developed since...
- * Recognising structure
 - * text → tree
 - * parse tree → AST
 - * forest disambiguation
- * tokens → list
- * image → visual model



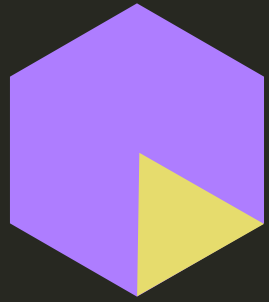
A.V. Aho & J.D. Ullman, The Theory of Parsing, Translation and Compiling, 1972.
V.Zaytsev, A.H.Bagge, Parsing in a Broad Sense, MoDELS 2014.



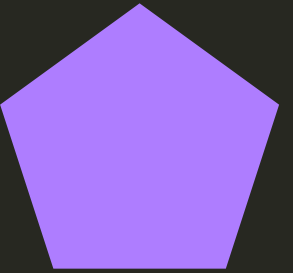
↑ Parsing



- * Reduce the input back to the start symbol
- * Recognise terminals
- * Replace terminals by nonterminals
- * Replace terminals and nonterminals by lhs
- * LR(1) ::= yacc | Beaver | Eli | SableCC | Irony;
- * GLR ::= bison | DMS | GDK | Tom;
- * SGLR ::= ASF+SDF | Spoofax | Stratego;

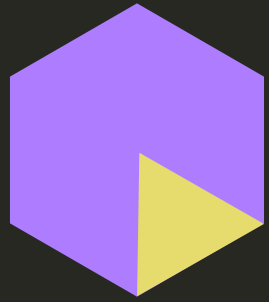


↓ Parsing

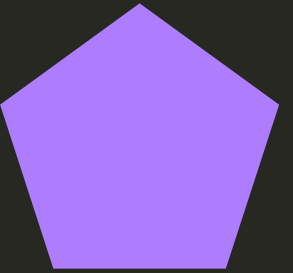


- * Imitate production by rederivation
- * Each nonterminal is a goal
- * Replace each goal by subgoals
- * Parse tree is built from top to bottom

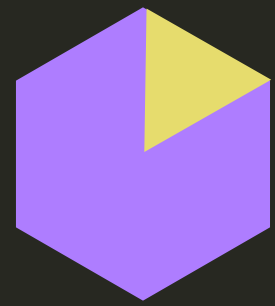
- * `LL(k) ::= JavaCC;` `LL(*) ::= ANTLR | TXL;`
- * `Earley ::= Marpa | ModelCC;` `DCG ::= Prolog;`
- * `GLL ::= Rascal | gll-combinators;`
- * `Packrat ::= Rats! | OMeta | PetitParser;`



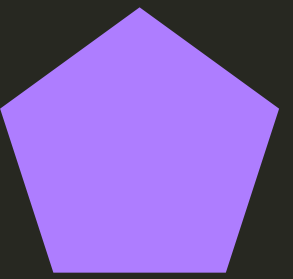
Semiparsing



- * grep
- * anchor terminals
- * islands & noise
- * skeleton grammars
- * relaxation & robustness
- * multilanguage

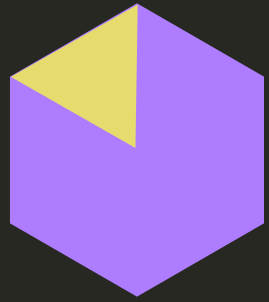


Fact extraction

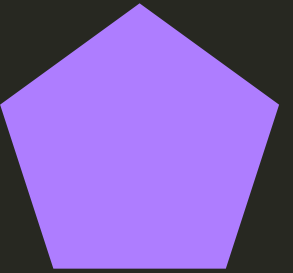


= parsing + generating a factbase
(or, sequence of graph transformations)

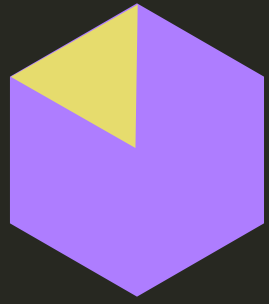
- * e.g., metrics
- * Can be language-parametric!
- * Schema
 - * describes form of the data
 - * ASG = Abstract Semantic Graph
 - * call graph
 - * dependence graph
 - * relations



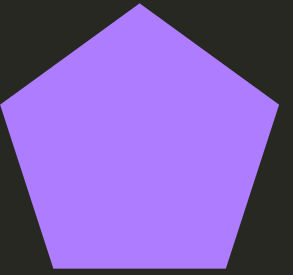
Slicing



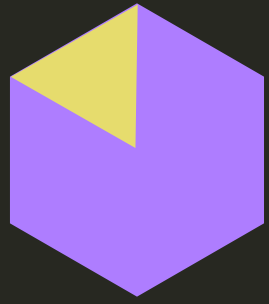
```
read(text);
read(n);
lines = 1;
chars = 1;
subtext = "";
c = getChar(text);
while (c != '\eof')
    if (c == '\n')
        then lines = lines + 1;
           chars = chars + 1;
    else chars = chars + 1;
        if (n != 0)
            then subtext = subtext ++ c;
               n = n - 1;
        c = getChar(text);
write(lines);
write(chars);
write(subtext);
```



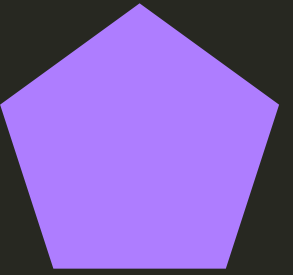
Slicing



- * Forward/backward slicing
- * Dynamic/conditioned slicing
 - * constraints on input
- * Chopping
 - * discover connection between I & O
- * Amorphous slicing
- * . . .



Slicing



- * Debugging

- * cf. Weiser CACM 1982

- * Cohesion measurement

- * cf. Ott&Bieman IST 1998

- * Comprehension

- * cf. De Lucia&Fasolino&Munro IWPC 1996

- * Maintenance

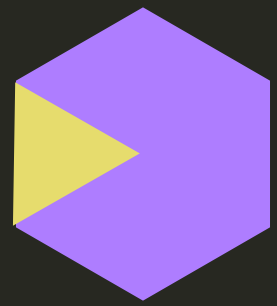
- * e.g. reuse

- * Re-engineering

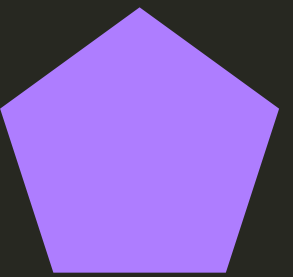
- * e.g. clone detection

Pattern matching

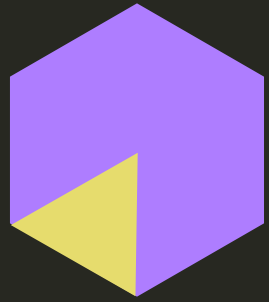
- * Easy to formulate on ADTs
- * In Rascal:
 - * `visit(){case}`
 - * `:=` and `!:=`
 - * functions
- * Need traversal strategies
 - * depth-first (pre-, in-, post-order)
 - * breadth-first
 - * topdown, bottomup, downup
 - * innermost, outermost
 - * . . .



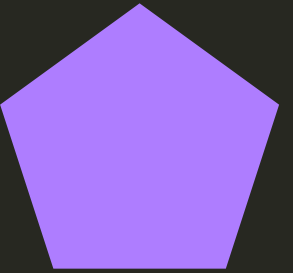
Decomposition



- * Recall partitioning & equiv. classes
- * Simplest form: modularisation
- * Usually: some graph + SCCs
- * Given granularity
 - * make a valid decomposition
 - * maximising benefit
- * Applicable to packages, build targets, automata, tasks, formulae, processes, rels...



Exploration



Software visualisation

Algorithm visualisation

Static
algorithm
visualisation

Algorithm
animation

Program visualisation

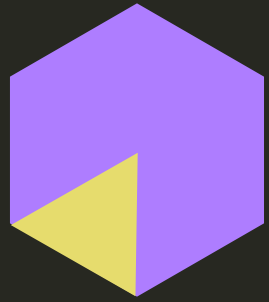
Data
animation

Static code
visualisation

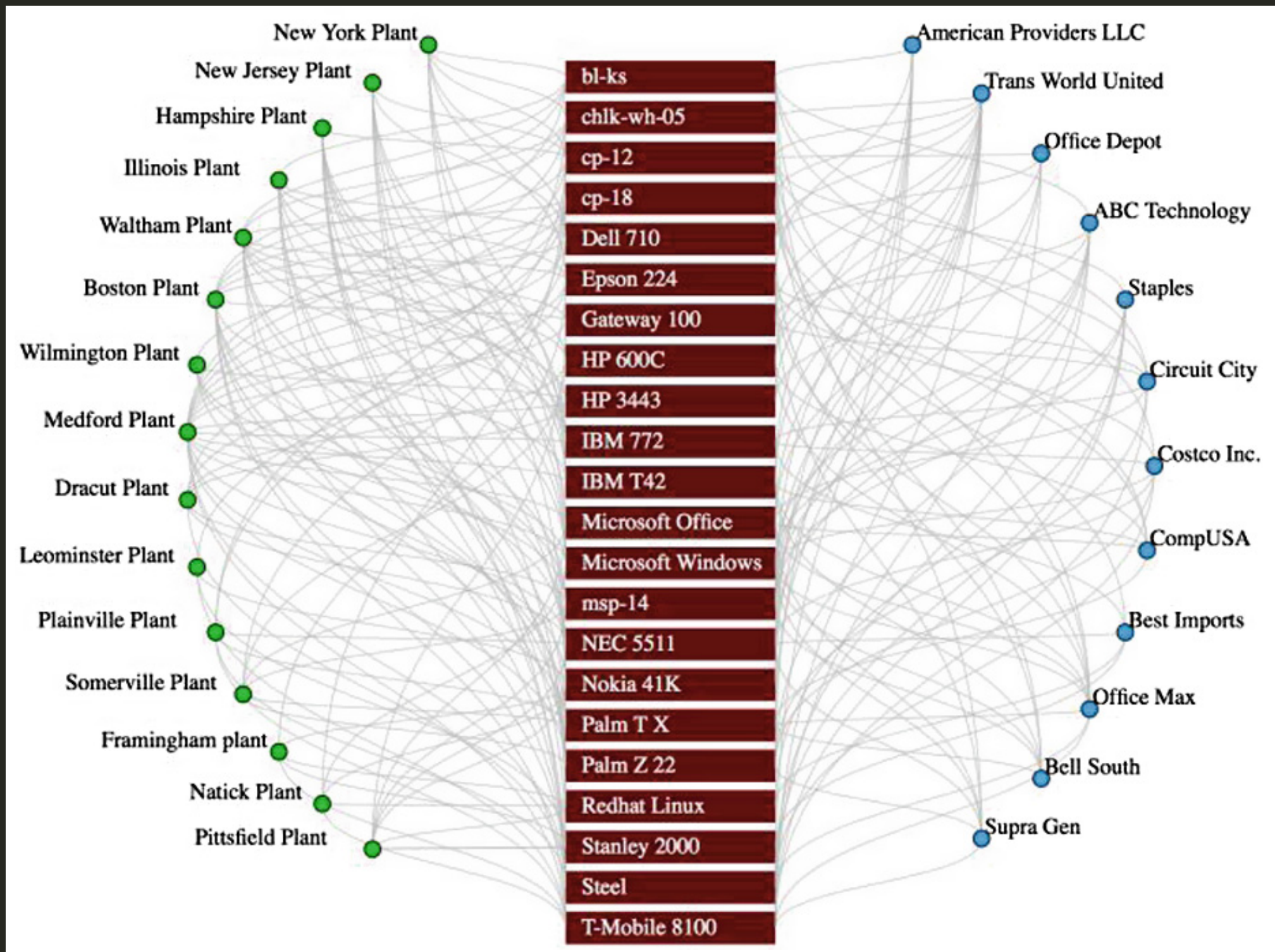
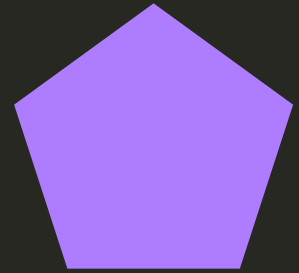
Static data
visualisation

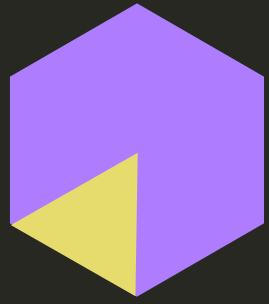
Visual programming

Code
animation

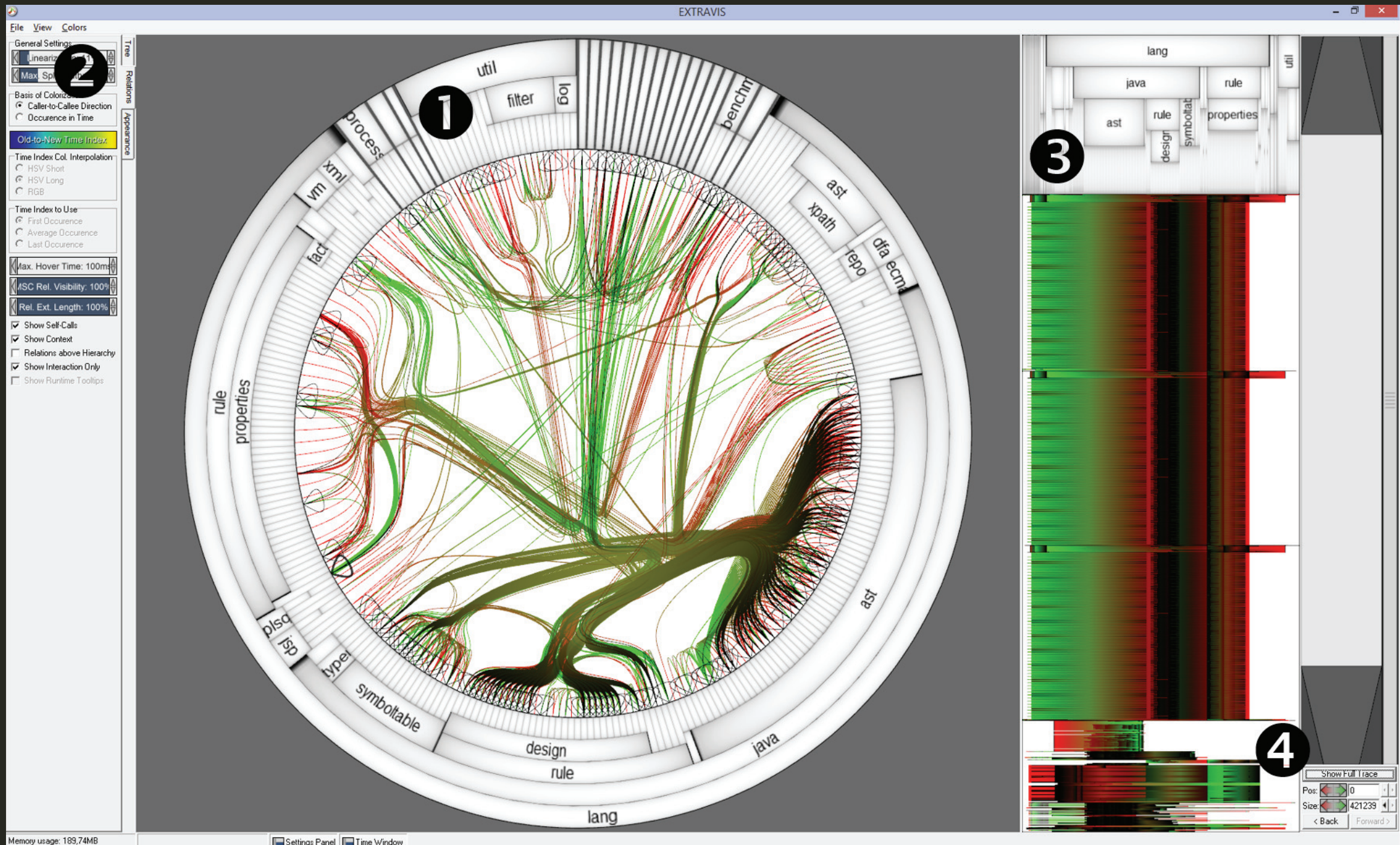
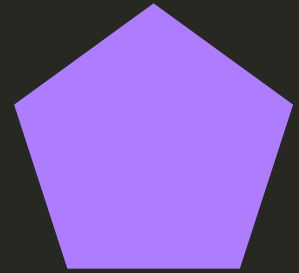


Visualisation

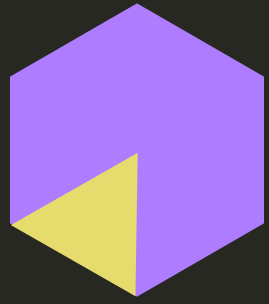




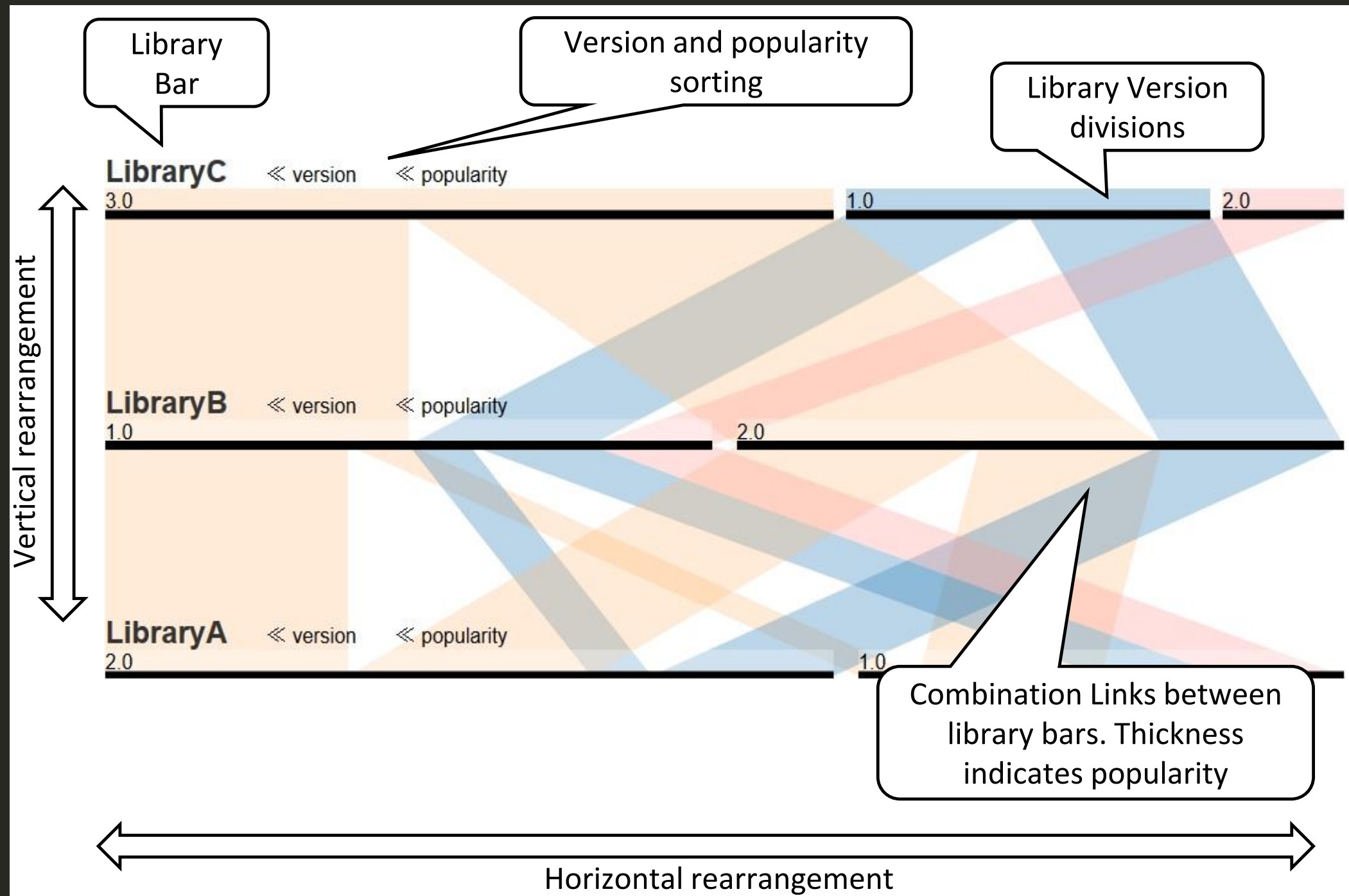
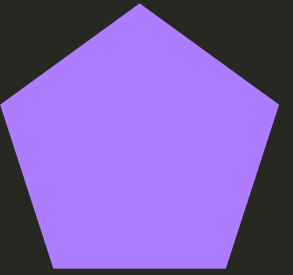
Trace vis

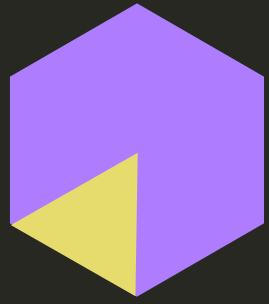


F.Fittkau, S.Finke, W.Hasselbring, J.Waller, Comparing trace visualizations for program comprehension through controlled experiments, ICPC 2015. <http://bibtex.github.io/ICPC-2015-FittkauFW.html>

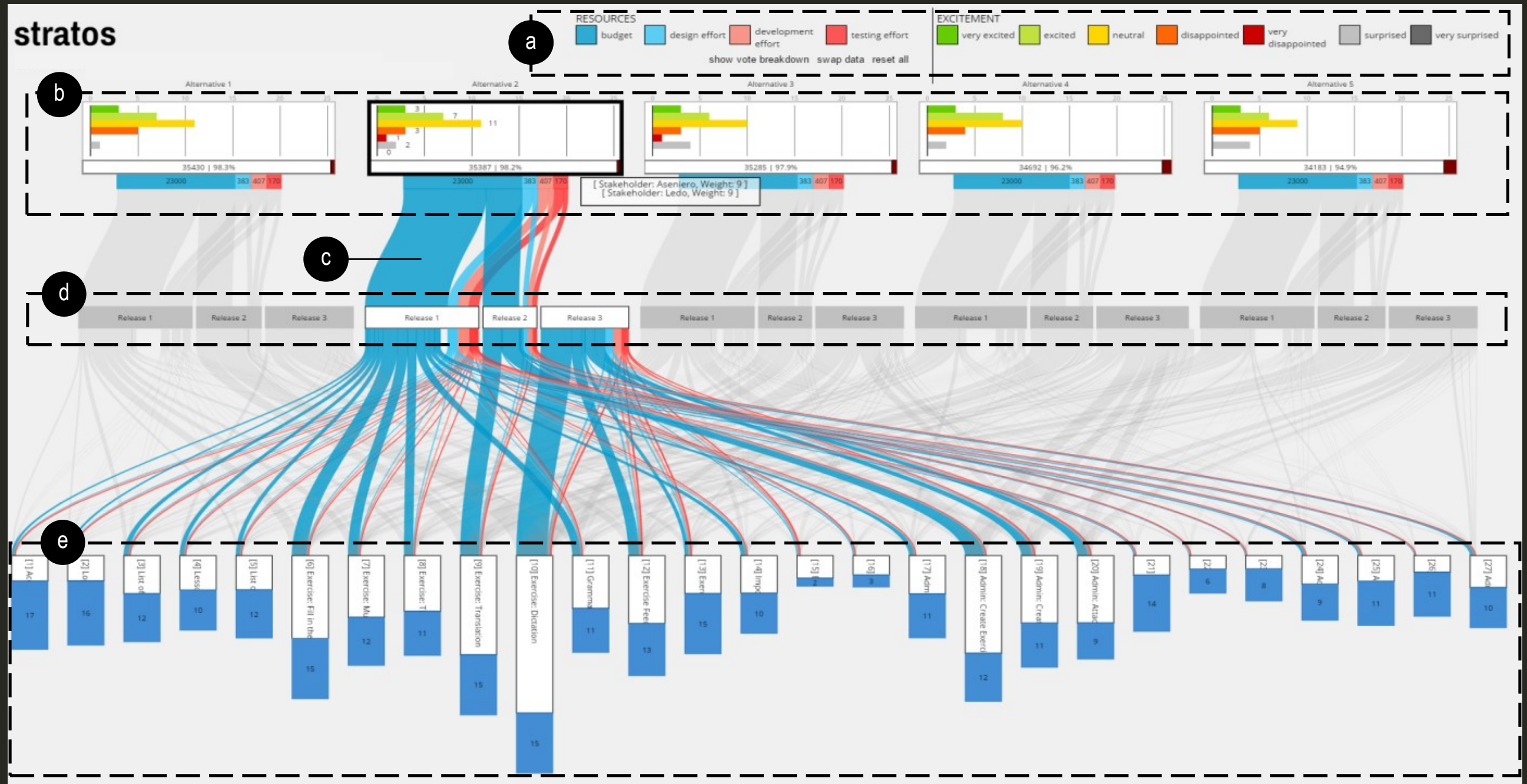
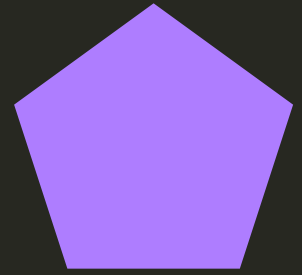


Versioning vis

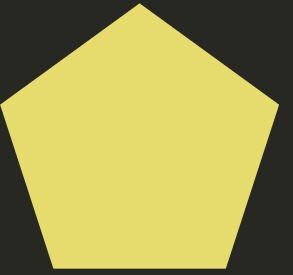
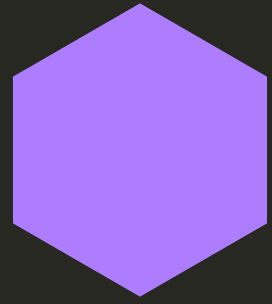




Release vis

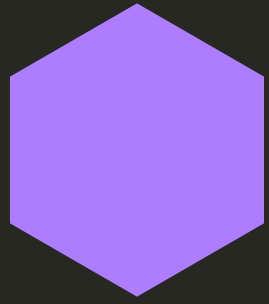


Data Reverse Engineering

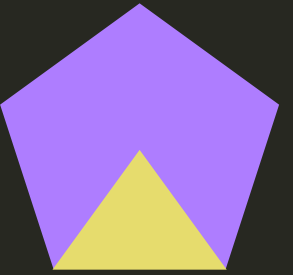


Data reverse engineering

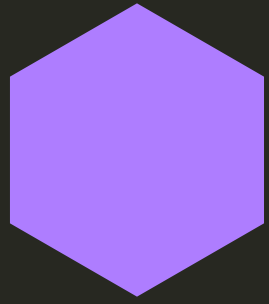
- * Database design recovery
- * Pattern recognition
- * Information retrieval
- * Clustering
- * Mining unstructured data



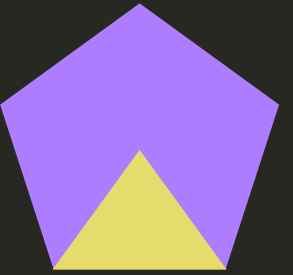
Database design recovery



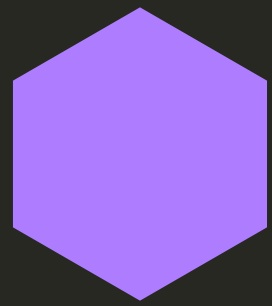
- * Forward database engineering
 - * Conceptual design
 - * Logical design
 - * Simplification
 - * Optimisation
 - * Translation
 - * Physical design
 - * View design



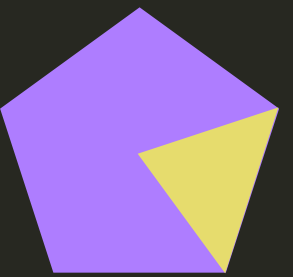
Database design recovery



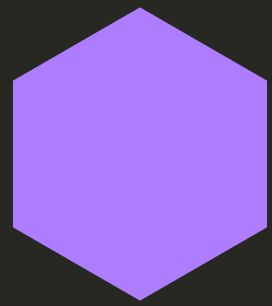
- * Data structure extraction
 - * Program analysis
 - * Data analysis
 - * Schema integration
- * Data structure conceptualisation
 - * Untranslation
 - * Deoptimisation
 - * Conceptual normalisation



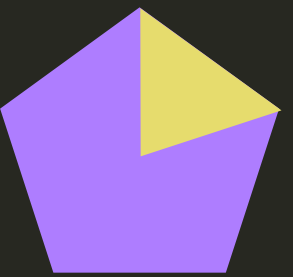
Pattern recognition



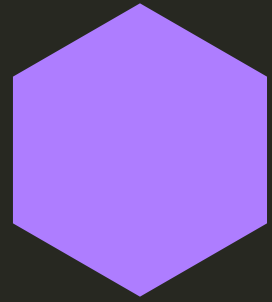
- * Pattern = feature vector
- * Quantitative features
 - * continuous / discrete / interval
- * Qualitative features
 - * nominal / ordinal
- * Find most descriptive/discriminatory



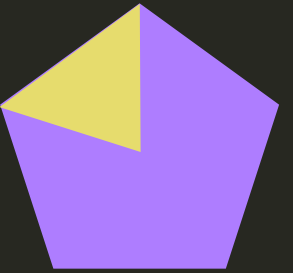
Information retrieval



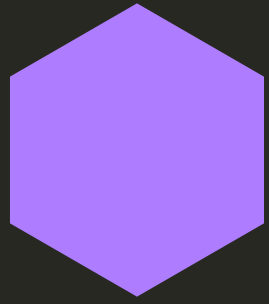
- * Knowledge discovery
- * Data mining
- * Usually statistical methods
 - * = require training
- * WEKA = Waikato Environment for Knowledge Analysis
 - * Java, 1992–2015
 - * <http://www.cs.waikato.ac.nz/ml/weka/>
 - * good with Groovy, Scala, Jython...



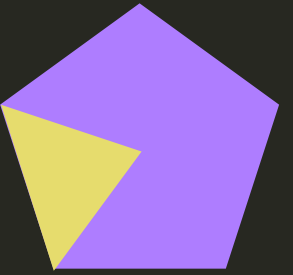
Clustering



- * Pattern recognition & representation
 - * similarity/proximity measure
 - * Minkowski / edit / statistical
- * Clustering techniques
 - * hierarchical / partitional
 - * agglomerative / divisive
 - * hard / fuzzy
 - * incremental / non-incremental
- * Dendrograms



MUD



- * Mixture

- * natural language text

- * technical artefacts

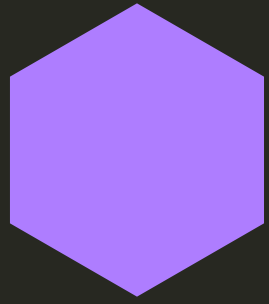
- * Unstructured data

- * dev communication

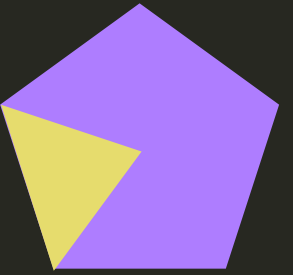
- * issue reports

- * documentation

- * meeting notes



MUD



- * Can fish for
 - * code fragments
 - * class names
 - * stack traces
 - * patches
 - * jargon
- * State of the art
 - * heuristic-based idiosyncratic tools

Conclusion

- * Besides forward engineering
 - * there is reverse engineering
- * Software comprehension
- * Code reverse engineering
 - * parsing, slicing, matching, visualising
- * Data reverse engineering
 - * design recovery, PR, IR, clustering, MUD
- * Mature yet active field