

Method and tool support for classifying software languages with Wikipedia

Ralf Lämmel, Dominik Mosen and Andrei Varanovich
Software Languages Team, University of Koblenz-Landau

<http://softlang.uni-koblenz.de/wikitax/>

Why and how to classify software languages?

planet-sl.org/sle2013/

The term "software language" refers to artificial languages used in software development. These include general-purpose programming languages, domain-specific languages, modeling and metamodeling languages, data models and ontologies. Examples include general purpose modeling languages such as SysML and UML, metamodeling frameworks such as Ecore, MOF or GOPRR, domain-specific modeling languages for business process modeling, such as BPMN, or embedded systems, such as Simulink or Modelica, and specialized XML-based and OWL-based languages and vocabularies. The term "software language" is intentionally broad; besides the above categories and examples, it also encompasses implicit approaches to language definition, such as APIs and collections of design patterns.

planet-sl.org/slebok/

SL(E) BOK 2.0

[Group](#) [Announcements](#) [Discussions](#) [Photos](#) [Videos](#) [Members](#)

SL(E)BOK 2.0

- [Home](#)
- [SL\(E\)](#)
- [BOK](#)
- [2.0](#)

Welcome

SL(E)BOK 2.0 is an emerging community-based collaborative-oriented **project** that aims at creating and maintaining a **Body Of Knowledge(BOK)** about **Software Languages, Software Linguistics and Software Language Engineering (SL(E))**.

Agenda

September 2012

[SL\(E\)BOK @ SLE2012](#) takes place on September 25th, 2012, at Dresden, Germany

101 companies: the emerging hitchhiker's guide through the software galaxy

http://101companies.org/wiki/Software_language

- ▶ Language:Parsec *instanceOf* this
- ◀ this *instanceOf* Namespace:Concept
- ◀ this *instanceOf* Vocabulary:Software language engineering
- ▶ Data manipulation language *isA* this
- ▶ Query language *isA* this
- ▶ Style sheet language *isA* this
- ▶ Tool-defined language *isA* this
- ▶ Transformation language *isA* this
- ▶ XML language *isA* this

BTW, where is the Wikipedia
for “Software Language”?

:-)

Anyone?

Problem-specific *exploration* tools

- Baskaya et al.: A tool for ontology-editing and **ontology**-based information exploration. ESAIR 2010.
- Haun et al.: CET: A tool for creative exploration of **graphs**. ECML/PKDD 2010.
- Dumas et al.: ViDaX: an interactive **semantic data** visualisation and exploration tool. AVI 2012.
- Hora et al.: Bug Maps: A tool for the visual exploration and analysis of **bugs**. CSMR 2012.
- De Roover et al.: Multi-dimensional exploration of **API** usage. ICPC 2013.

Category graph exploration with ***WikiTax***

<http://softlang.uni-koblenz.de/wikitalx/>

Category graph exploration with WikiTax

The screenshot displays the WikiTax application window, which is used for exploring category graphs. The interface is divided into several sections:

- File View:** Located at the top left, it contains a tree view of categories. The 'Computer languages' category is expanded, showing subcategories like 'Lists of computer languages', 'Articles with example code', 'Data modeling languages', 'XML', 'Data types', 'Markup languages', 'Programming language topics', 'Programming languages', 'Specification languages', 'Algorithm description languages', 'Data modeling languages', 'XML', 'Dependently typed languages', 'Formal specification languages', 'Hardware description languages', 'Unified Modeling Language', 'Stylesheet languages', and 'Transformation languages'. The 'Articles with example code' category is highlighted in blue.
- Categories overall:** A summary of the overall category statistics, showing 132 categories and 591 pages.
- Pages overall:** A summary of the overall page statistics, showing 591 pages.
- Selected node:** The current selected node is 'Articles with example code', which has 44 pages, 29 subcategories, 44 pages (transitive), 29 subcategories (transitive), and 1 parent category.
- Highlighting:** A legend for the highlighting mode, showing color-coded boxes for different page counts: 0 pages (grey), ≤ 14 pages (white), ≤ 45 pages (light red), ≤ 66 pages (red), and ≤ 204 pages (dark red).
- Parent Categories:** A list of parent categories, currently showing 'Computer languages'.
- Pages:** A list of pages associated with the selected node, including 'A* search algorithm', 'ARM architecture', 'Babbage (programming language)', 'Batch file', and 'Binary decision diagram'.
- Highlighting mode:** A dropdown menu at the bottom left, currently set to 'pages'.
- remove marked categories:** A button at the bottom right to remove marked categories.

Computer languages



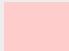
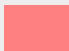

- ▶ ☒ Lists of computer languages
- ▶ ☒ Articles with example code
- ▼ ☐ Data modeling languages
 - ☒ XML
- ▶ ☒ Data types
- ▶ ☐ Markup languages
- ▶ ☒ Programming language topics
- ▶ ☐ Programming languages
- ▼ ☐ Specification languages
 - ☐ Algorithm description languages
 - ▼ ☐ Data modeling languages
 - ☒ XML
 - ☐ Dependently typed languages
 - ☐ Formal specification languages
 - ☐ Hardware description languages
 - ☒ Unified Modeling Language
- ▶ ☐ Stylesheet languages
- ▶ ☐ Transformation languages

Result of 2 levels of *extraction* with
some categories marked for *exclusion*

Categories overall:	132
Pages overall:	591
Selected node:	Articles w
pages:	44
subcategories:	29
pages (transitive):	44
subcategories (transitive):	29
parent categories:	1
comment:	

Edit comment ...

Highlighting

	,... = 0 pages
	≤ 14 pages
	≤ 45 pages
	≤ 66 pages
	≤ 204 pages

Parent Categories

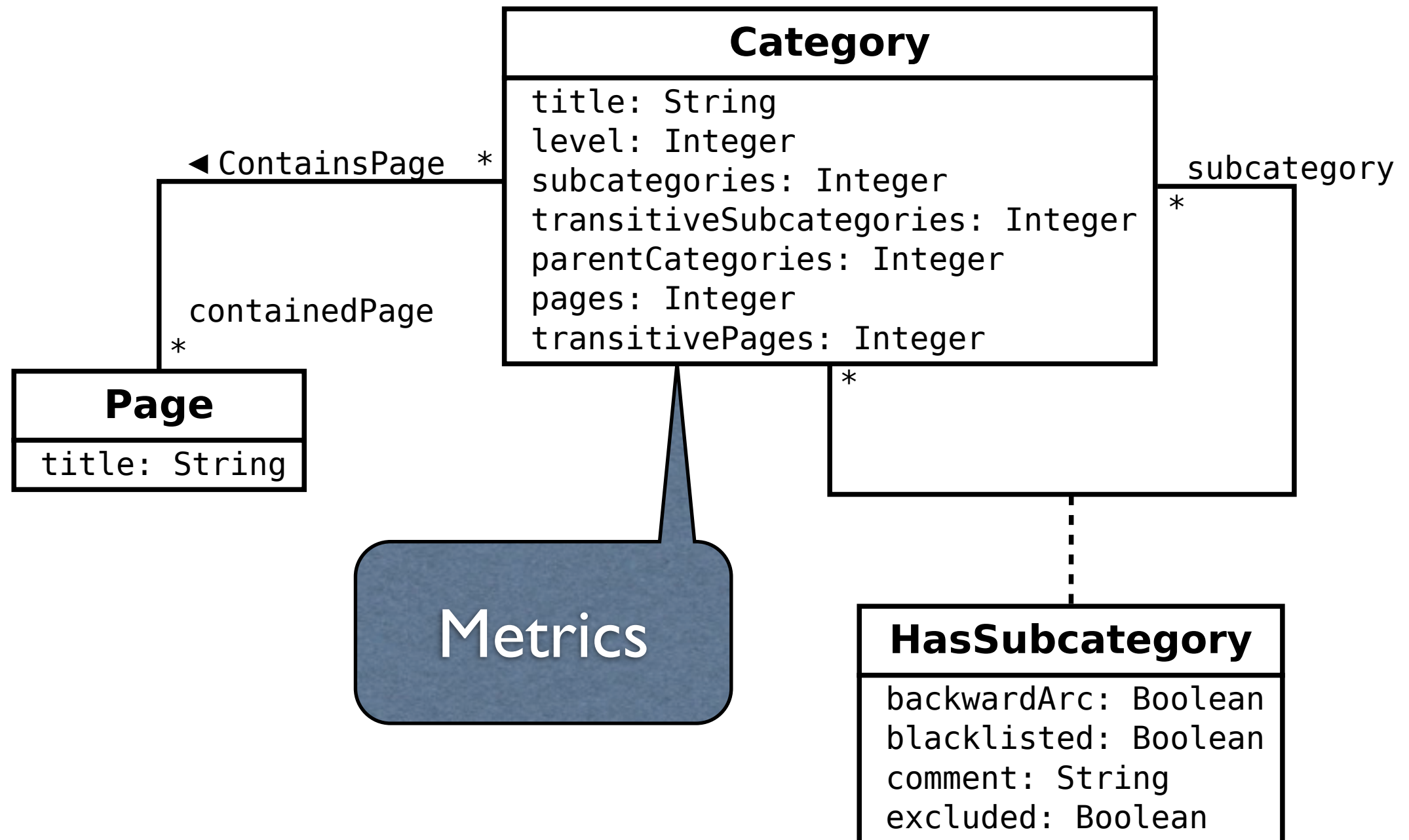
Computer languages

Pages

The *WikiTax* approach

- Category graph **extraction** (from *Wikipedia*)
- ... **reduction** (by the exclusion of categories)
- ... **visualization** (using simple metrics)
- ... **export** (for external processing)

Metamodel of the *WikiTax* category graph



Implementation of *WikiTax*

- Wikipedia API
- TGraphs, JSON, and CSV
- Java (Swing)
- JGraLab

See GitHub URL etc.:
<http://softlang.uni-koblenz.de/wikitax/>
(Open Source and Open Data)

Problem-specific
concern: ***exclusion***

Computer languages

- ▶ ☒ Lists of computer languages
- ▶ ☒ Articles with example code
- ▼ ☐ Data modeling languages
 - ☒ XML
- ▶ ☒ Data types
- ▶ ☐ Markup languages
- ▶ ☒ Programming language topics
- ▶ ☐ Programming languages
- ▼ ☐ Specification languages
 - ☐ Algorithm description languages
 - ▼ ☐ Data modeling languages
 - ☒ XML
 - ☐ Dependently typed languages
 - ☐ Formal specification languages
 - ☐ Hardware description languages
 - ☒ Unified Modeling Language
- ▶ ☐ Stylesheet languages
- ▶ ☐ Transformation languages

Why and what and
how to exclude?

Are these
true classifiers
in terms of
software concepts?

Exclusion types

- Alternative classifier (unrelated to software concepts)
 - ▶ e.g., Academic programming languages
- Deviating classifier (in fact, non-classifier)
 - ▶ e.g., Articles with example code
- Singleton classifier (focusing on one language)
 - ▶ e.g., Cascading Style Sheets
- List classifier (collecting list pages)
 - ▶ e.g., Lists of programming languages
- Maintenance classifier
 - ▶ e.g., Uncategorized programming languages

Category	Exclusion type
<i>Academic programming languages</i>	Alternative classifier
<i>Articles with example code</i>	Deviating classifier
<i>Cascading Style Sheets</i>	Singleton classifier
<i>Data types</i>	Deviating classifier
<i>Discontinued programming languages</i>	Alternative classifier
<i>DocBook</i>	Singleton classifier
<i>Esoteric programming languages</i>	Alternative classifier
<i>Experimental programming languages</i>	Alternative classifier
<i>HTML</i>	Singleton classifier
<i>JSON</i>	Singleton classifier
<i>Lists of computer languages</i>	List classifier
<i>Lists of programming languages</i>	List classifier
<i>Markup language comparisons</i>	Deviating classifier
<i>Markup language stubs</i>	Maintenance classifier
<i>Non-English-based programming languages</i>	Alternative classifier
<i>Programming language families</i>	Deviating classifier
<i>Programming language standards</i>	Deviating classifier
<i>Programming language topics</i>	Deviating classifier
<i>Programming languages by creation date</i>	Alternative classifier
<i>Programming languages conferences</i>	Deviating classifier
<i>Software by programming language</i>	Deviating classifier
<i>SyncML</i>	Singleton classifier
<i>TeX</i>	Singleton classifier
<i>Text Encoding Initiative</i>	Singleton classifier
<i>Troff</i>	Singleton classifier
<i>Uncategorized programming languages</i>	Maintenance classifier
<i>Unified Modeling Language</i>	Singleton classifier
<i>Wikipedia categories named after programming languages</i>	Deviating classifier
<i>XML</i>	Singleton classifier

Computer languages

Category	Subcategories
<i>Data modeling languages</i>	–
<i>Markup languages</i>	<i>Declarative markup languages, GIS file formats, Knowledge representation languages, Lightweight markup languages, Mathematical markup languages, Musical markup languages, Page description markup languages, Playlist markup languages, User interface markup languages, Vector graphics markup languages, Web syndication formats, XML markup languages</i>
<i>Programming languages</i>	<i>.NET programming languages, Agent-based programming languages, Agent-oriented programming languages, Concatenative programming languages, Concurrent programming languages, Data-structured programming languages, Declarative programming languages, Dependently typed languages, Domain-specific programming languages, Dynamic programming languages, Extensible syntax programming languages, Formula manipulation languages, Function-level languages, Functional languages, High Integrity Programming Language, High-level programming languages, ICL programming languages, Intensional programming languages, Low-level programming languages, Multi-paradigm programming languages, Nondeterministic programming languages, Object-based programming languages, Pattern matching programming languages, Procedural programming languages, Process termination functions, Prototype-based programming languages, Reactive programming languages, Secure programming languages, Set theoretic programming languages, Statically typed programming languages, Synchronous programming languages, Term-rewriting programming languages, Text-oriented programming languages, Tree programming languages, Visual programming languages, XML-based programming languages</i>
<i>Specification languages</i>	<i>Algorithm description languages, Dependently typed languages, Formal specification languages, Hardware description languages</i>
<i>Stylesheet languages</i>	–
<i>Transformation languages</i>	<i>Macro programming languages</i>

Category	Subcategories
<i>Data modeling languages</i>	–
<i>Markup languages</i>	<i>Declarative markup languages, GIS file formats, Knowledge representation languages, Lightweight markup languages, Mathematical markup languages, Musical markup languages, Page description markup languages, Playlist markup languages, User interface markup languages, Vector graphics markup languages, Web syndication formats, XML markup languages</i>
<i>Programming languages</i>	<i>.NET programming languages, Agent-based programming languages, Agent-oriented programming languages, Concatenative programming languages, Concurrent programming languages, Data-structured programming languages, Declarative programming languages, Dependently typed languages, Domain-specific programming languages, Dynamic programming languages, Extensible syntax programming languages, Formula manipulation languages, Function-level languages, Functional languages, High Integrity Programming Language, High-level programming languages, ICL programming languages, Intensional programming languages, Low-level programming languages, Multi-paradigm programming languages, Nondeterministic programming languages, Object-based programming languages, Pattern matching programming languages, Procedural programming languages, Process termination functions, Prototype-based programming languages, Reactive programming languages, Secure program-</i>

An aside:
Where are the
query languages?

native programming languages, Concurrent programming languages, Data-structured programming languages, Declarative programming languages, Dependently typed languages, Domain-specific programming languages, Dynamic programming languages, Extensible syntax programming languages, Formula manipulation languages, Function-level languages, Functional languages, High Integrity Programming Language, High-level programming languages, ICL programming languages, Intensional programming languages, Low-level programming languages, Multi-paradigm programming languages, Nondeterministic programming languages, Object-based programming languages, Pattern matching programming languages, Procedural programming languages, Process termination functions, Prototype-based programming languages, Reactive programming languages, Secure programming languages, Set theoretic programming languages, Statically typed programming languages, Synchronous programming languages, Term-rewriting programming languages, Text-oriented programming languages, Tree programming languages, Visual programming languages, XML-based programming languages

Specification languages

Algorithm description languages, Dependently typed languages, Formal specification languages, Hardware description languages

Stylesheet languages

—

Transformation languages

Macro programming languages

Let's compare
SQL and XSLT!

<http://en.wikipedia.org/wiki/SQL>

Categories: Database management systems Computer languages Data modeling languages Declarative programming languages **Query languages** Relational database management systems SQL

Category: Query languages

Categories: Domain-specific programming languages
Data management Databases

A level-2
subcategory of
computer
languages

<http://en.wikipedia.org/wiki/XSLT>

Categories: Declarative programming languages Functional languages Markup languages Transformation languages World Wide Web Consortium standards XML-based programming languages XML-based standards

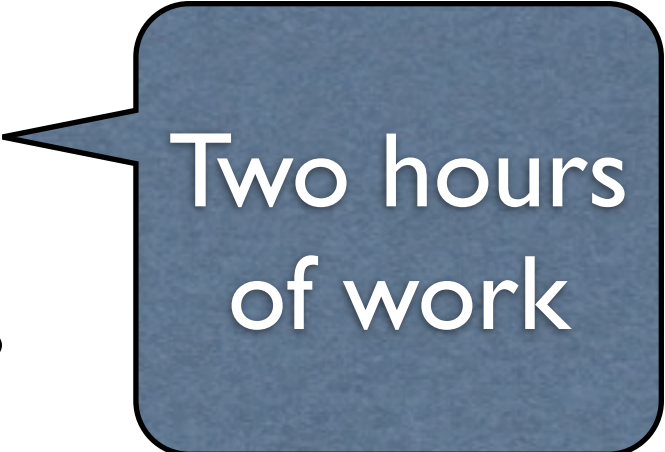
A level-1
subcategory of
computer
languages

Category: Transformation language

Categories: Computer languages

Programming languages -- all levels

- Initial extraction
 - ▶ 423 categories, 7515 pages, 8 levels
- 1st pruning phase
 - ▶ 29 excluded categories as discussed earlier
 - ▶ 288 categories, 6671 pages
- 2nd pruning phase
 - ▶ 79 categories, 1560 pages, 4 levels



Two hours
of work

Pages

Programming languages

- ☐ ICL programming languages
- ☐ Agent-based programming languages
- ☐ Agent-oriented programming languages
- ☐ Concatenative programming languages
- ☐ Concurrent programming languages
- ▶ ☐ Data-structured programming languages
- ▶ ☐ Declarative programming languages
- ☐ Dependently typed languages
- ▶ ☐ Domain-specific programming languages
- ▶ ☐ Dynamic programming languages
- ☐ Extensible syntax programming languages
- ☐ Formula manipulation languages
- ☐ Function-level languages
- ☐ Functional languages
- ☐ High Integrity Programming Language
- ☐ High-level programming languages
- ☐ Intensional programming languages
- ☐ Low-level programming languages
- ▶ ☐ Multi-paradigm programming languages
- ☐ .NET programming languages
- ☐ Nondeterministic programming languages
- ▶ ☐ Object-based programming languages
- ☐ Pattern matching programming languages
- ▶ ☐ Procedural programming languages
- ☐ Process termination functions
- ☐ Prototype-based programming languages
- ☐ Reactive programming languages
- ☐ Secure programming languages
- ☐ Set theoretic programming languages
- ☐ Statically typed programming languages
- ☐ Synchronous programming languages
- ☐ Term-rewriting programming languages
- ☐ Text-oriented programming languages
- ☐ Tree programming languages
- ☐ Visual programming languages
- ☐ XML-based programming languages

Categories

Programming languages

- ☐ ICL programming languages
- ☐ Agent-based programming languages
- ☐ Agent-oriented programming languages
- ☐ Concatenative programming languages
- ☐ Concurrent programming languages
- ▶ ☐ Data-structured programming languages
- ▶ ☐ Declarative programming languages
- ☐ Dependently typed languages
- ▶ ☐ Domain-specific programming languages
- ▶ ☐ Dynamic programming languages
- ☐ Extensible syntax programming languages
- ☐ Formula manipulation languages
- ☐ Function-level languages
- ☐ Functional languages
- ☐ High Integrity Programming Language
- ☐ High-level programming languages
- ☐ Intensional programming languages
- ☐ Low-level programming languages
- ▶ ☐ Multi-paradigm programming languages
- ☐ .NET programming languages
- ☐ Nondeterministic programming languages
- ▶ ☐ Object-based programming languages
- ☐ Pattern matching programming languages
- ▶ ☐ Procedural programming languages
- ☐ Process termination functions
- ☐ Prototype-based programming languages
- ☐ Reactive programming languages
- ☐ Secure programming languages
- ☐ Set theoretic programming languages
- ☐ Statically typed programming languages
- ☐ Synchronous programming languages
- ☐ Term-rewriting programming languages
- ☐ Text-oriented programming languages
- ☐ Tree programming languages
- ☐ Visual programming languages
- ☐ XML-based programming languages

Programming languages

- ☐ ICL programming languages
- ☐ Agent-based programming languages
- ☐ Agent-oriented programming languages
- ☐ Concatenative programming languages
- ☐ Concurrent programming languages
- ▶ ☐ Data-structured programming languages
- ▶ ☐ Declarative programming languages
- ☐ Dependently typed languages
- ▶ ☐ Domain-specific programming languages
- ▶ ☐ Dynamic programming languages
- ☐ Extensible syntax programming languages
- ☐ Formula manipulation languages
- ☐ Function-level languages
- ☐ Functional languages
- ☐ High Integrity Programming Language
- ☐ High-level programming languages
- ☐ Intensional programming languages
- ☐ Low-level programming languages
- ▶ ☐ Multi-paradigm programming languages
- ☐ .NET programming languages
- ☐ Nondeterministic programming languages
- ▶ ☐ Object-based programming languages
- ☐ Pattern matching programming languages
- ▶ ☐ Procedural programming languages
- ☐ Process termination functions
- ☐ Prototype-based programming languages
- ☐ Reactive programming languages
- ☐ Secure programming languages

Programming languages

- ☐ ICL programming languages
- ☐ Agent-based programming languages
- ☐ Agent-oriented programming languages
- ☐ Concatenative programming languages
- ☐ Concurrent programming languages
- ▶ ☐ Data-structured programming languages
- ▶ ☐ Declarative programming languages
- ☐ Dependently typed languages
- ▶ ☐ Domain-specific programming languages
- ▶ ☐ Dynamic programming languages
- ☐ Extensible syntax programming languages
- ☐ Formula manipulation languages
- ☐ Function-level languages
- ☐ Functional languages
- ☐ High Integrity Programming Language
- ☐ High-level programming languages
- ☐ Intensional programming languages
- ☐ Low-level programming languages
- ▶ ☐ Multi-paradigm programming languages
- ☐ .NET programming languages
- ☐ Nondeterministic programming languages
- ▶ ☐ Object-based programming languages
- ☐ Pattern matching programming languages
- ▶ ☐ Procedural programming languages
- ☐ Process termination functions
- ☐ Prototype-based programming languages
- ☐ Reactive programming languages
- ☐ Secure programming languages

- ▶ ☐ Declarative programming languages
- ☐ Dependently typed languages
- ▶ ☐ Domain-specific programming languages
- ▶ ☐ Dynamic programming languages
- ☐ Extensible syntax programming languages
- ☐ Formula manipulation languages
- ☐ Function-level languages
- ☐ Functional languages
- ☐ High Integrity Programming Language
- ☐ High-level programming languages
- ☐ Intensional programming languages
- ☐ Low-level programming languages
- ▶ ☐ Multi-paradigm programming languages
- ☐ .NET programming languages
- ☐ Nondeterministic programming languages
- ▶ ☐ Object-based programming languages
- ☐ Pattern matching programming languages
- ▶ ☐ Procedural programming languages
- ☐ Process termination functions
- ☐ Prototype-based programming languages
- ☐ Reactive programming languages
- ☐ Secure programming languages
- ☐ Set theoretic programming languages
- ☐ Statically typed programming languages
- ☐ Synchronous programming languages
- ☐ Term-rewriting programming languages
- ☐ Text-oriented programming languages
- ☐ Tree programming languages
- ☐ Visual programming languages
- ☐ XML-based programming languages

- ▶ ☐ Declarative programming languages
- ☐ Dependently typed languages
- ▶ ☐ Domain-specific programming languages
- ▶ ☐ Dynamic programming languages
- ☐ Extensible syntax programming languages
- ☐ Formula manipulation languages
- ☐ Function-level languages
- ☐ Functional languages
- ☐ High Integrity Programming Language
- ☐ High-level programming languages
- ☐ Intensional programming languages
- ☐ Low-level programming languages
- ▶ ☐ Multi-paradigm programming languages
- ☐ .NET programming languages
- ☐ Nondeterministic programming languages
- ▶ ☐ Object-based programming languages
- ☐ Pattern matching programming languages
- ▶ ☐ Procedural programming languages
- ☐ Process termination functions
- ☐ Prototype-based programming languages
- ☐ Reactive programming languages
- ☐ Secure programming languages
- ☐ Set theoretic programming languages
- ☐ Statically typed programming languages
- ☐ Synchronous programming languages
- ☐ Term-rewriting programming languages
- ☐ Text-oriented programming languages
- ☐ Tree programming languages
- ☐ Visual programming languages
- ☐ XML-based programming languages

Thanks! Questions?

<http://softlang.uni-koblenz.de/wikitax/>