

# Tracking the Flow of Ideas through the Programming Languages Literature

Michael Greenberg, Kathleen Fisher, and David Walker





How can we understand  
the PL literature?

Is there more **related work** should I cite?

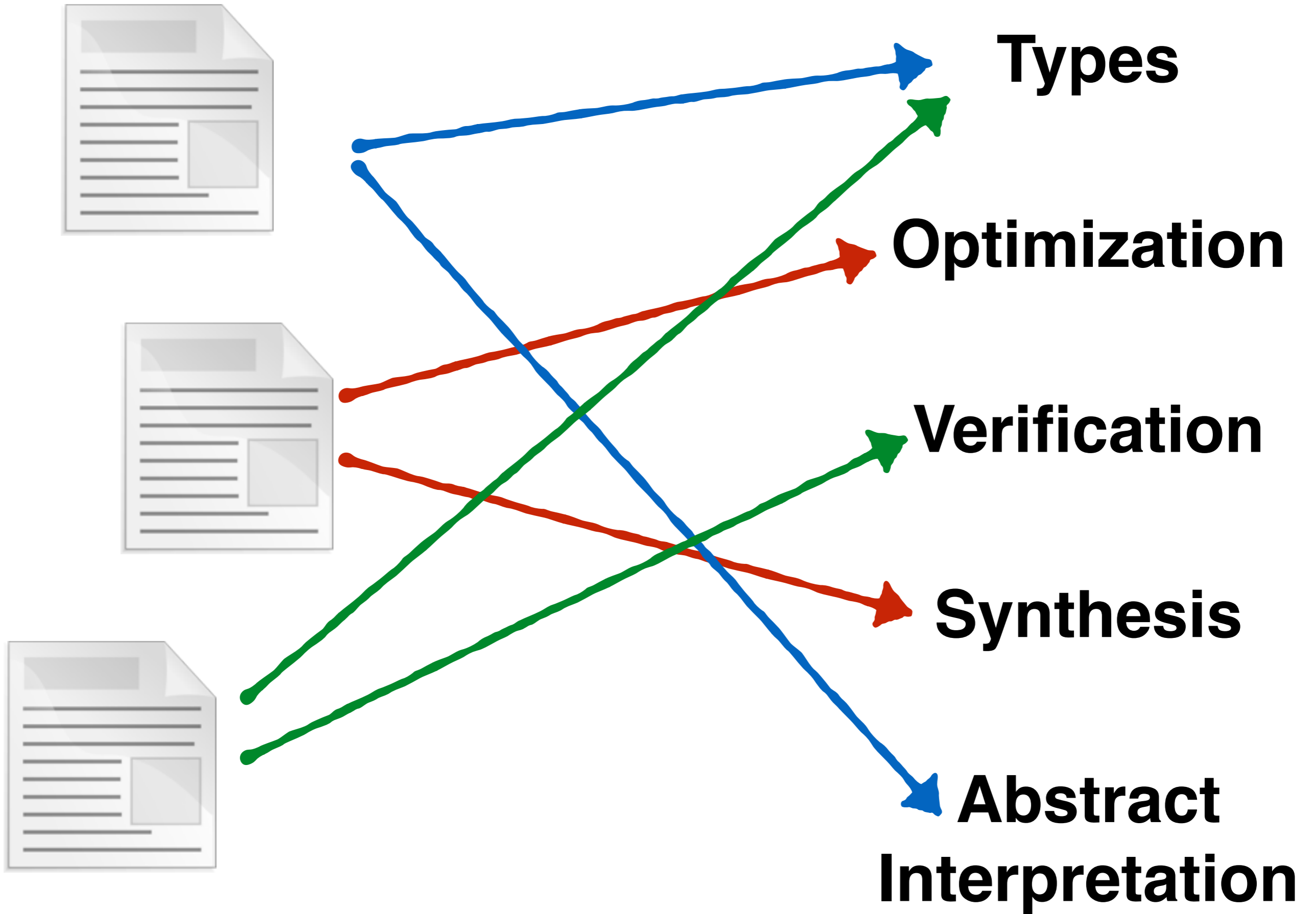
Is my work a better fit for **PLDI** or **POPL**?

Who should I invite to this **PC**?

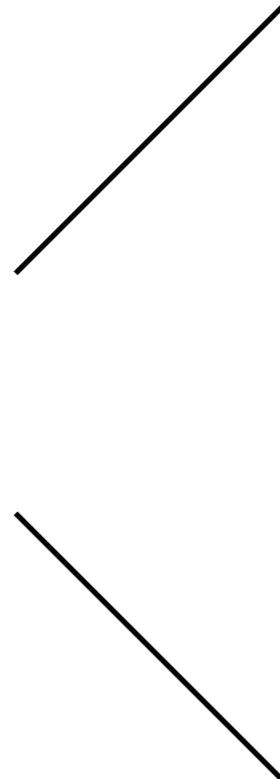
Who should **review** this paper?

Was this a **typical year** for ICFP?

How has **OOPSLA** changed over the years?

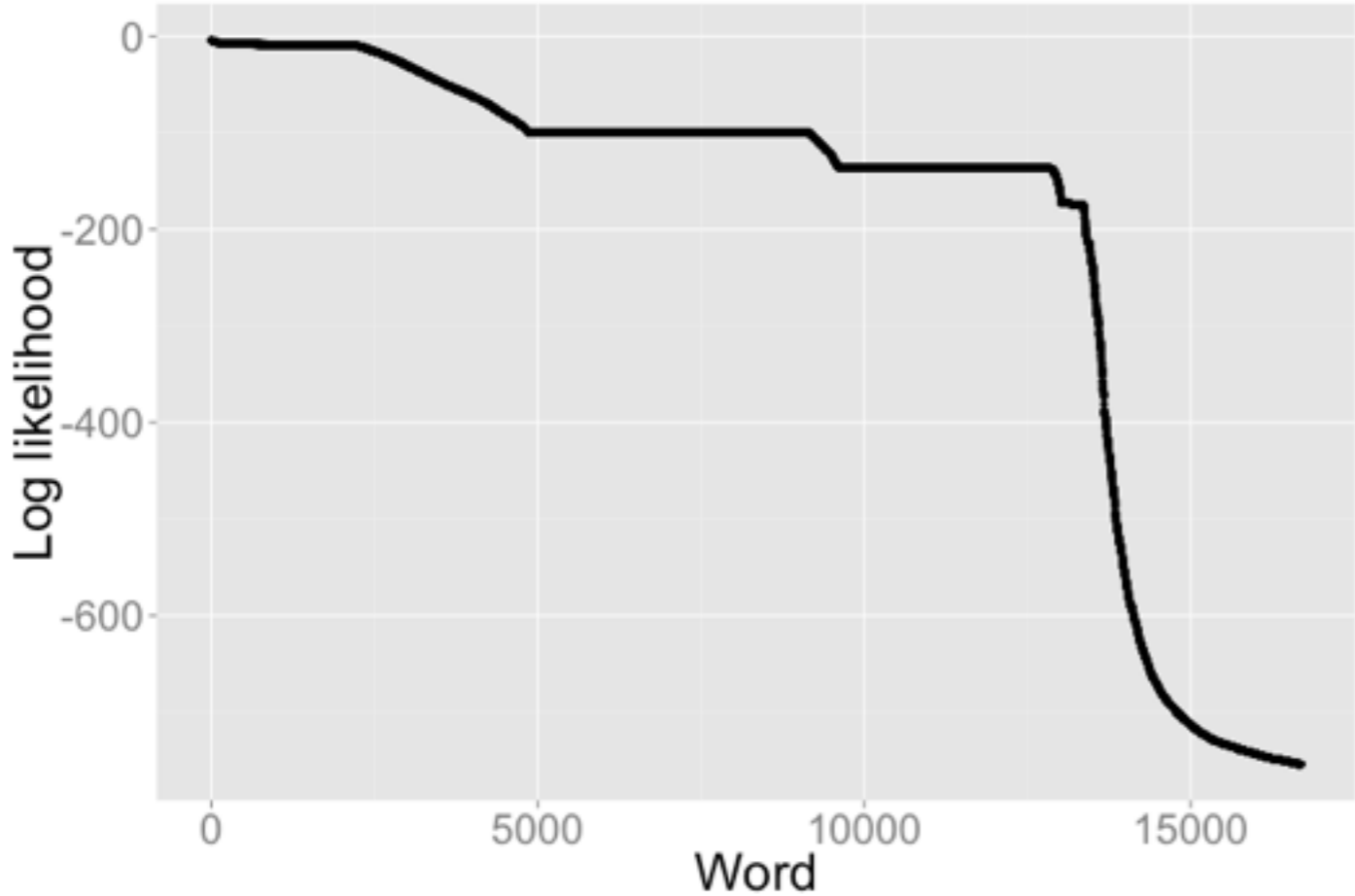


# What is a 'topic' in a document?



Word	Count
type	120
system	83
check	34
static	21

# Topics are distributions of words



"Parsing" topic	
Word	Log likelihood
grammar	-3.905040
language	-4.206531
structure	-4.308618
parser	-4.513348
...	...

# Documents are a mix of topics



**type systems** .6

Word	Count
type	120
system	83
check	34
static	21

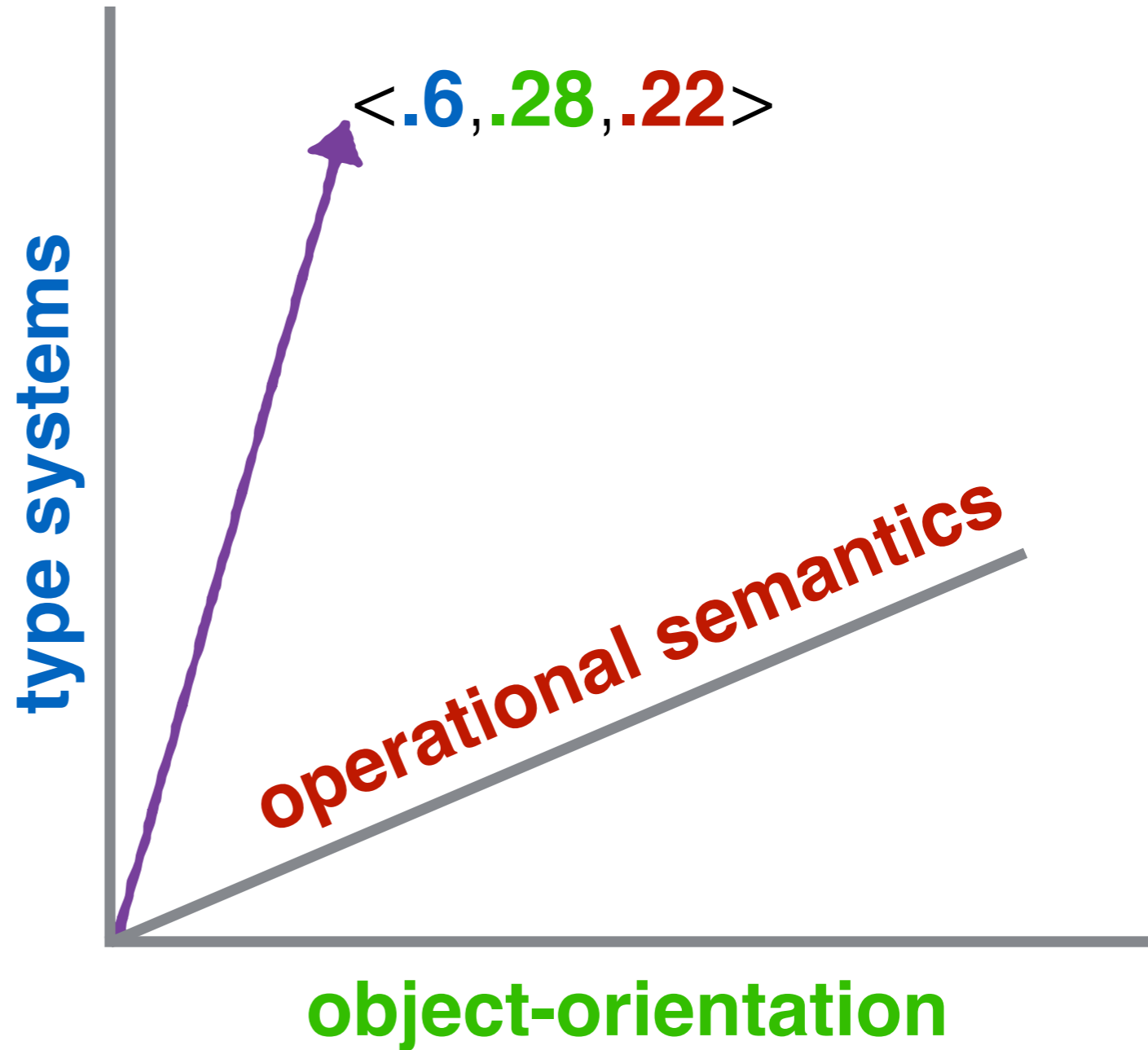
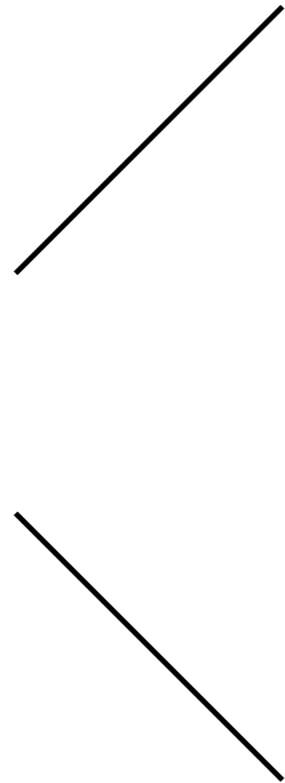
**object-orientation** .22

Word	Count
object	88
class	13
instance	12
method	7

**operational semantics** .28

Word	Count
semantics	90
step	45
reduce	38
evaluate	19

# Documents are a mix of topics





# Generative LDA topic model

## Gradual Typing for First-Class Classes\*

### Abstract

Dynamic type-checking and object-oriented programming often go hand-in-hand; scripting languages such as Python, Ruby, and JavaScript all embrace object-oriented (OO) programming. When scripts written in such languages grow and evolve into large programs, the lack of a static type discipline reduces maintainability. A programmer may thus wish to migrate parts of such scripts to a sister language with a static type system. Unfortunately, existing type systems neither support the flexible OO composition mechanisms found in scripting languages nor accommodate sound interoperation with untyped code.

In this paper, we present the design of a gradual typing system that supports sound interaction between statically- and dynamically-typed units of class-based code. The type system uses row polymorphism for classes and thus supports mixin-based OO composition. To protect migration of mixins from typed to untyped components, the system employs a novel form of contracts that partially seal classes. The design comes with a theorem that guarantees the soundness of the type system even in the presence of untyped components.



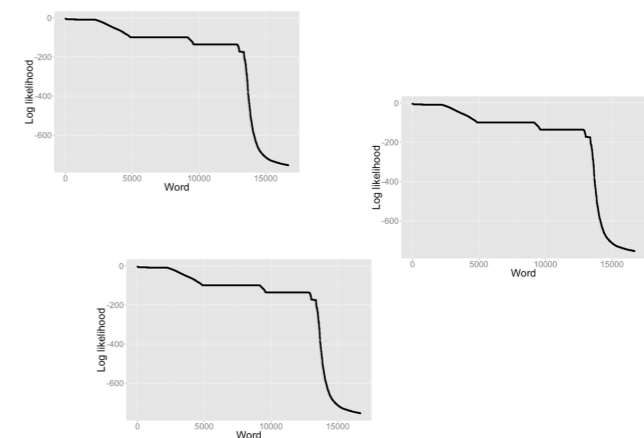
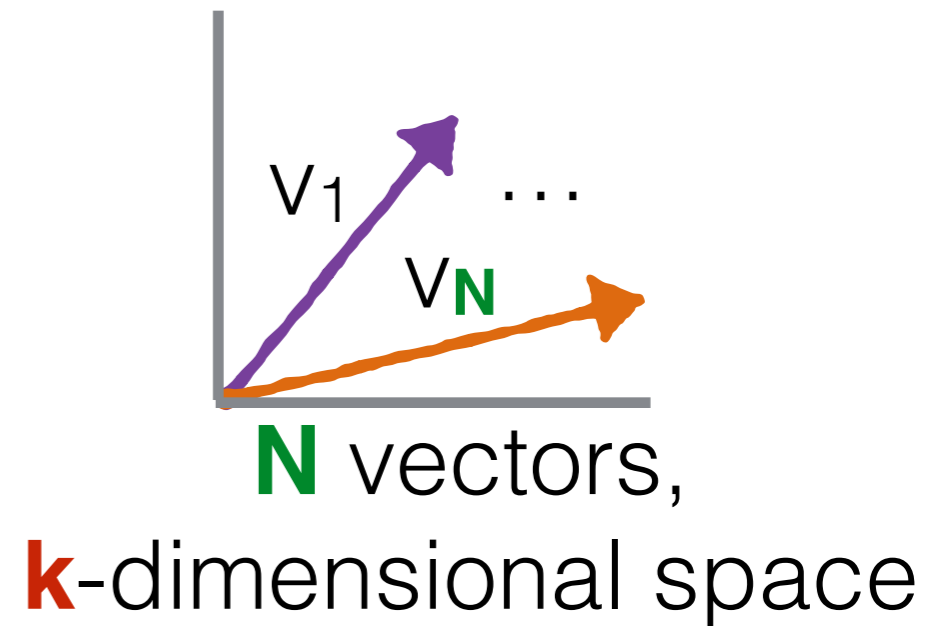
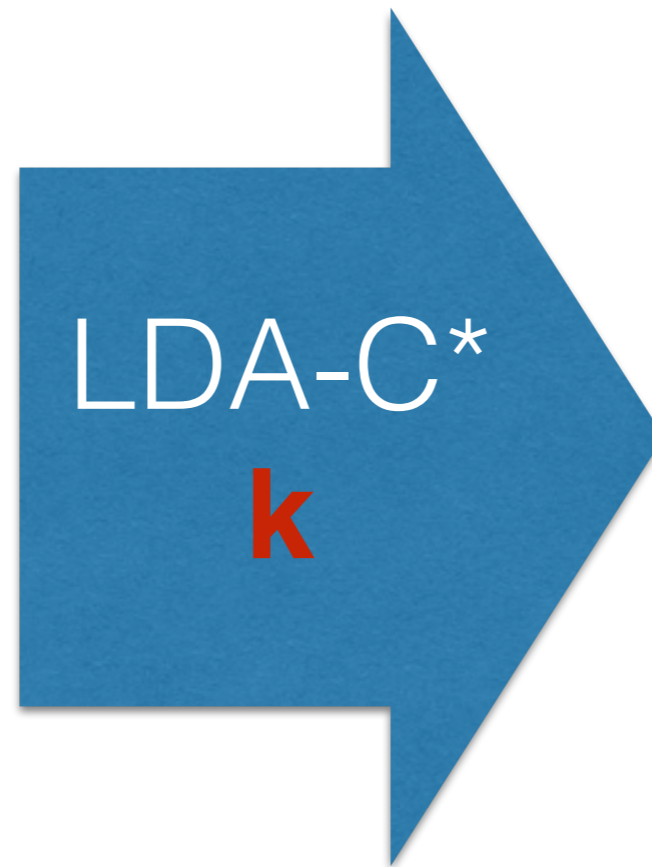
OOPSLA'12, October 19–26, 2012, Tucson, Arizona, USA.  
Copyright © 2012 ACM 978-1-4503-1561-6/12/10...\$10.00

Takikawa, Strickland, Dimoulas, Tobin-Hochstadt, and Felleisen  
*Gradual typing for first-class classes*. OOPSLA 2012.

# Inference with LDA



**N** bags of words



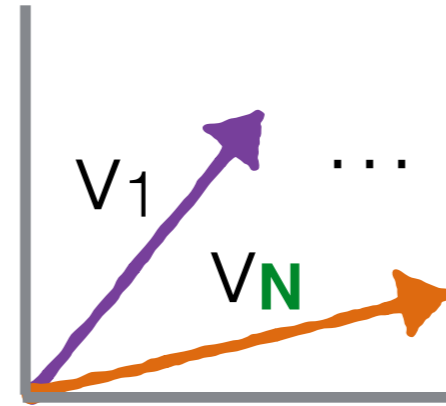
**k** topics



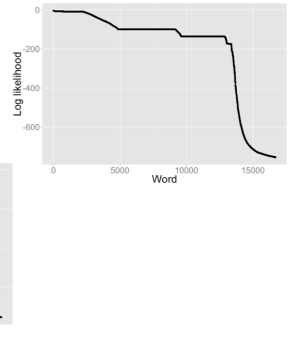
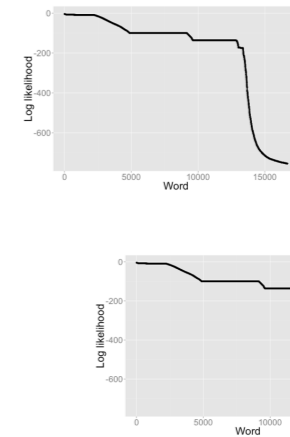
corpus  
**N** docs



**N** bags of words



**N** vectors



**k** topics



**k** top words    aggregate vectors  
**k** top papers    by year  
by conference



combined vocabulary



**k** topic names

# Parsing

- Parsing drops standard **stopwords**
  - Added some extra ones with TF-IDF
- **Stemmed** words using nltk\*
  - Removes plurals, etc.

a  
about  
above  
after  
again  
against  
...

calculi → calculus  
goes → go

\*<http://www.nltk.org/>

# Our corpora

- **Abstracts:** ICFP, OOPSLA, PLDI, POPL
  - 4,355 documents
  - Imperfect data in the ACM Digital Library
- **Fulltext:** PLDI, POPL
  - 2,257 documents
  - Imperfect PDF-to-text conversion

# Let's name a topic!

object  
Space overhead  
ent with  
heap  
region  
n  
memor  
r systems  
pointer  
collection  
collector  
garbage  
collection  
and their  
allocation  
and garbage collection gridlock  
reference  
on bounding time and space for multiprocessor garbage collection  
Garbage collection without paging

# Garbage collection!

# Topic names for k=20, abstracts

Compiler optimization

Array Processing

Verification

Program Logics

Resource management

Garbage Collection

Test generation

Parallelism

Parsing

Components and APIs

Object-Oriented Programming

Language Design

Low-level compiler optimizations

Program Analysis

Analysis of Concurrent Programs

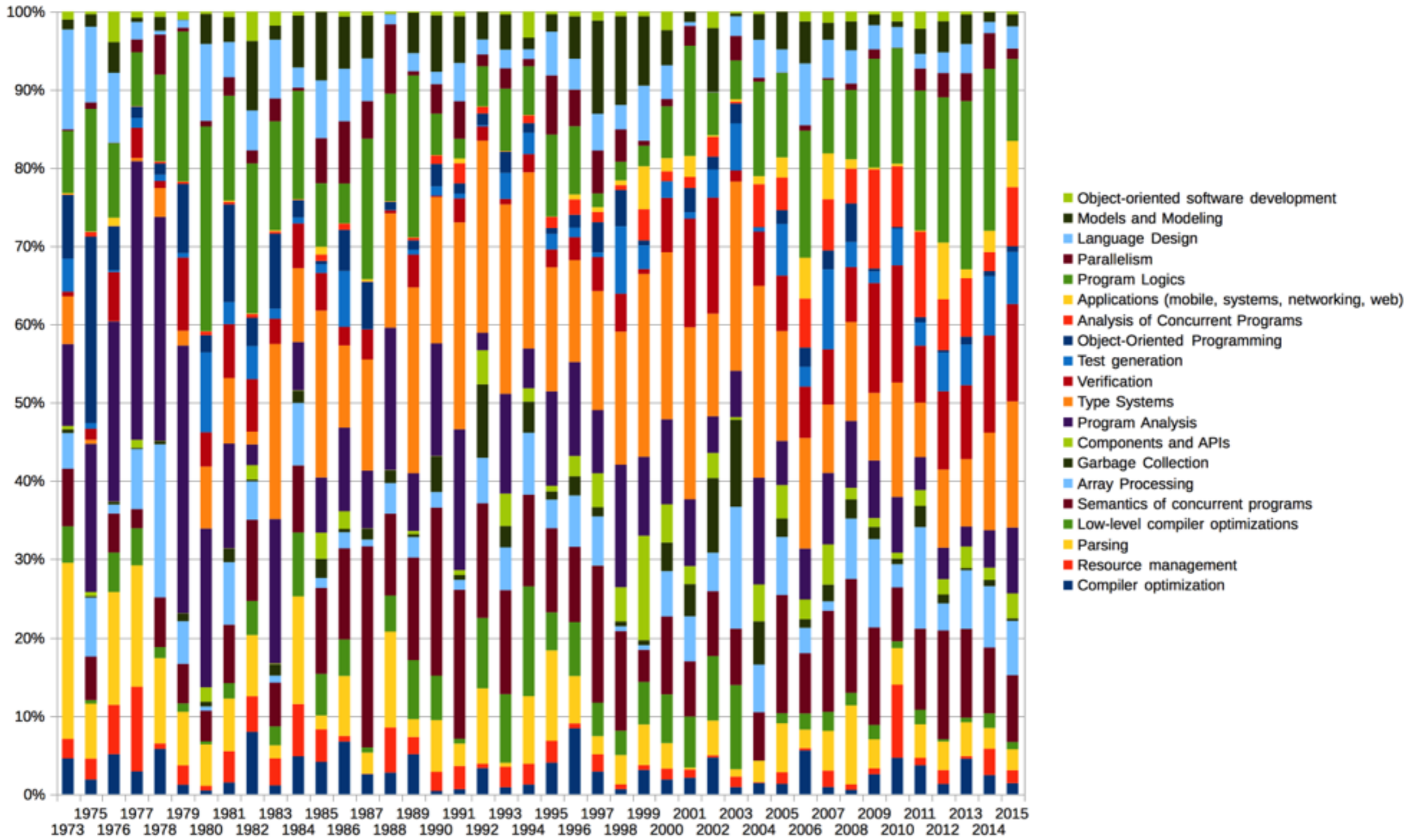
Models and Modeling

Semantics of concurrent programs

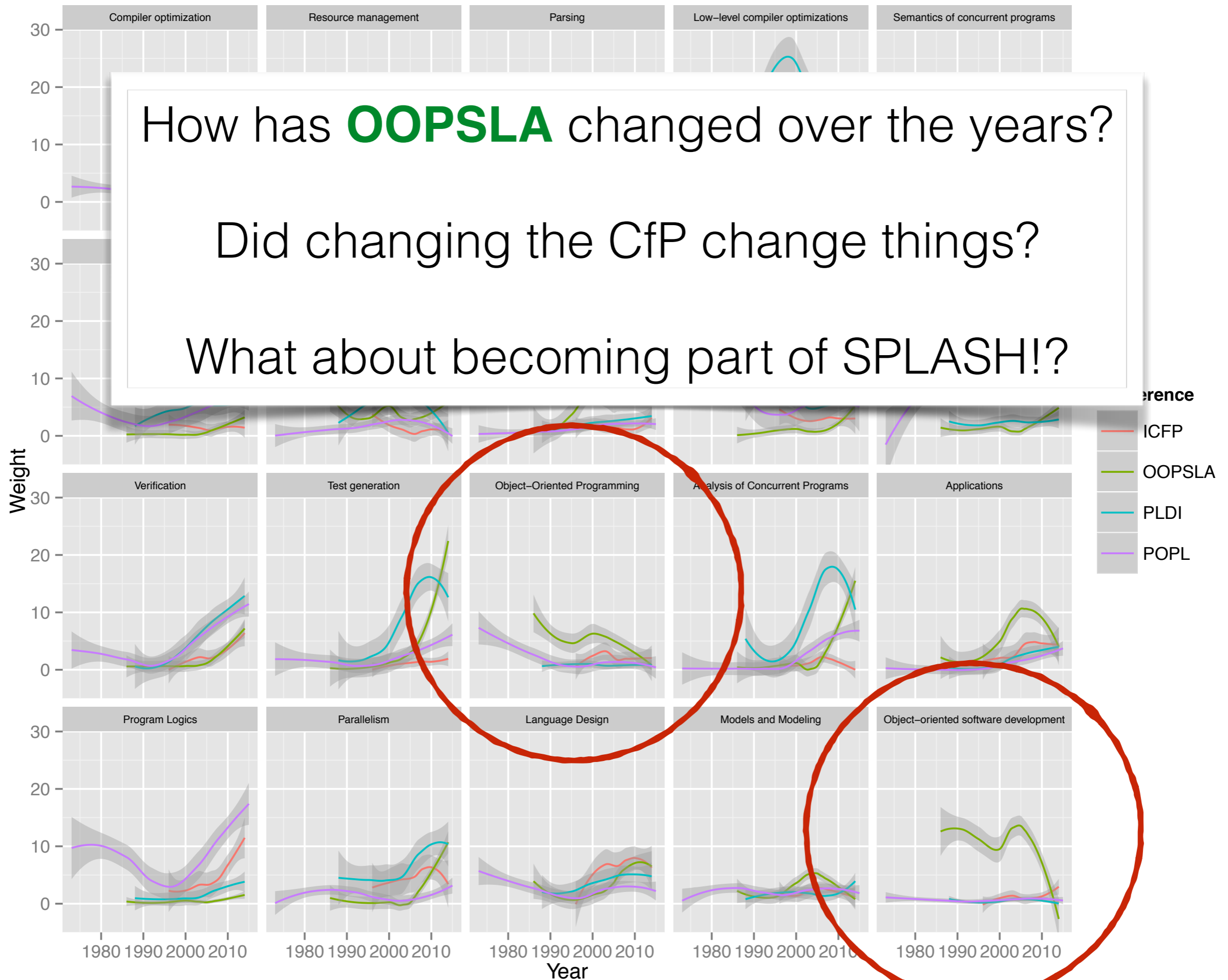
Type Systems

Applications

Object-oriented software development







# OOPSLA Call for Papers

2006

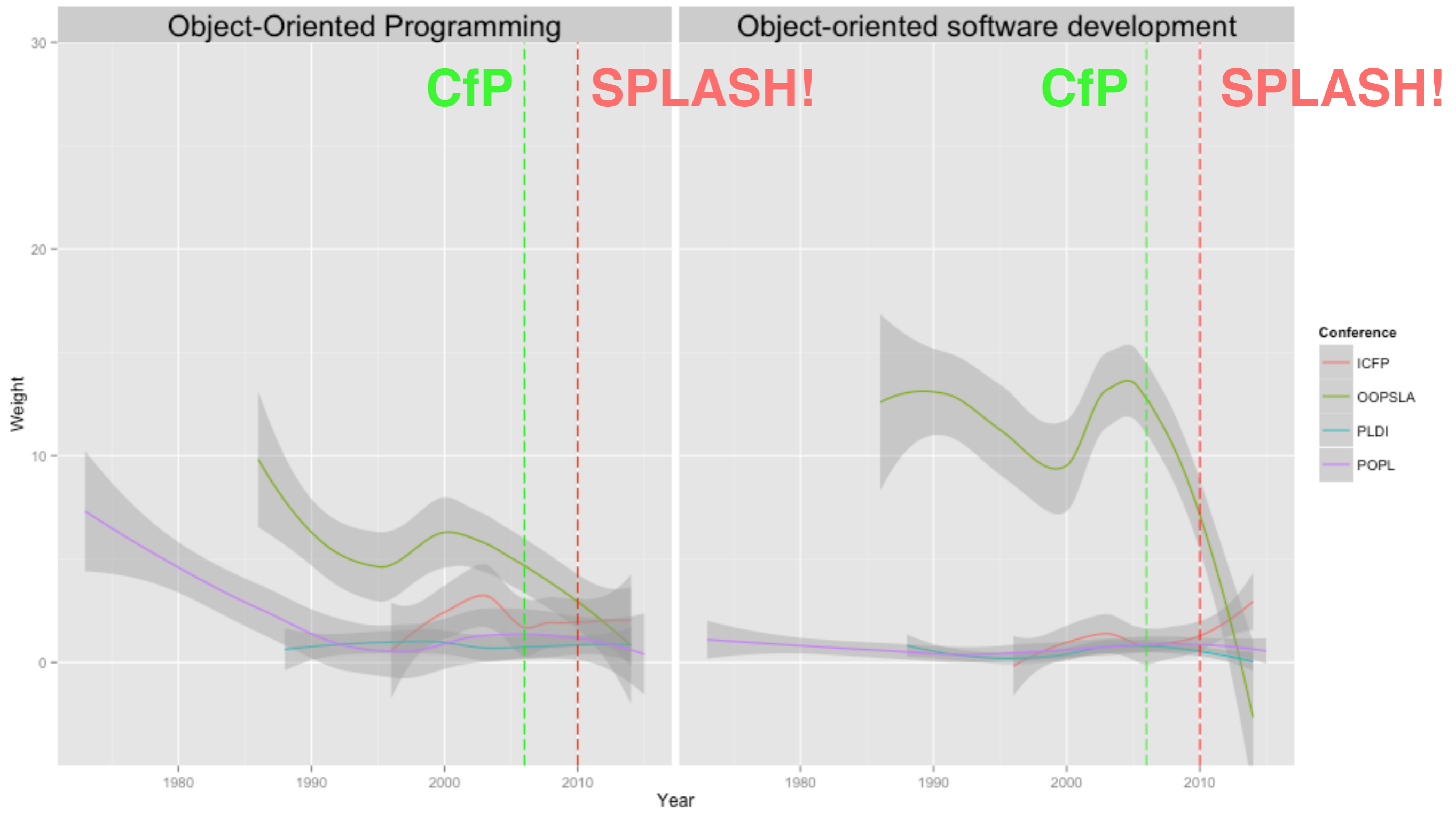
2007

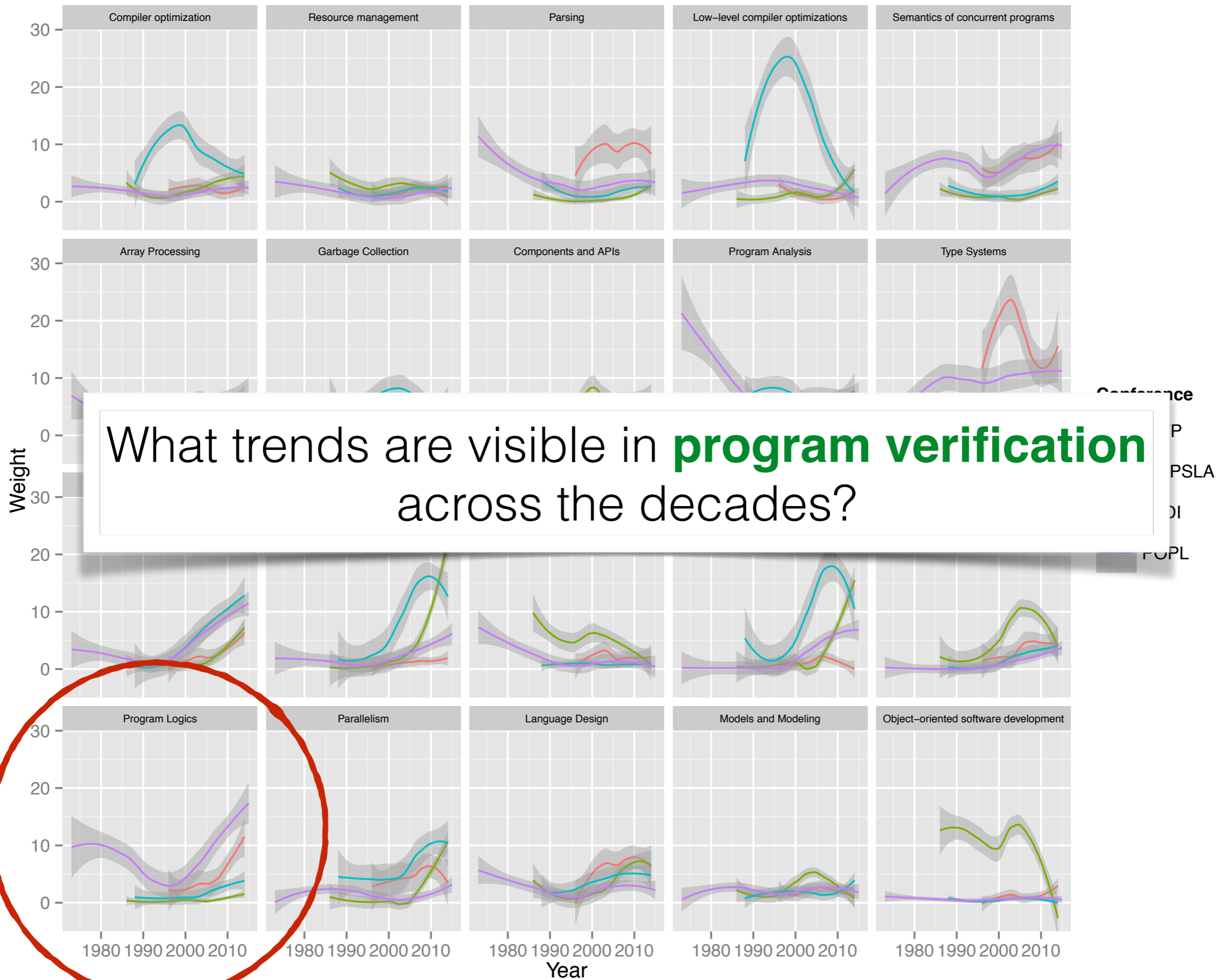
2010

foundations of **object  
and related  
technologies**

paradigms **beyond  
the traditional  
concept** of object-  
oriented programming

**all aspects** of  
programming  
languages and  
software engineering,  
**broadly construed**





# Program Logics

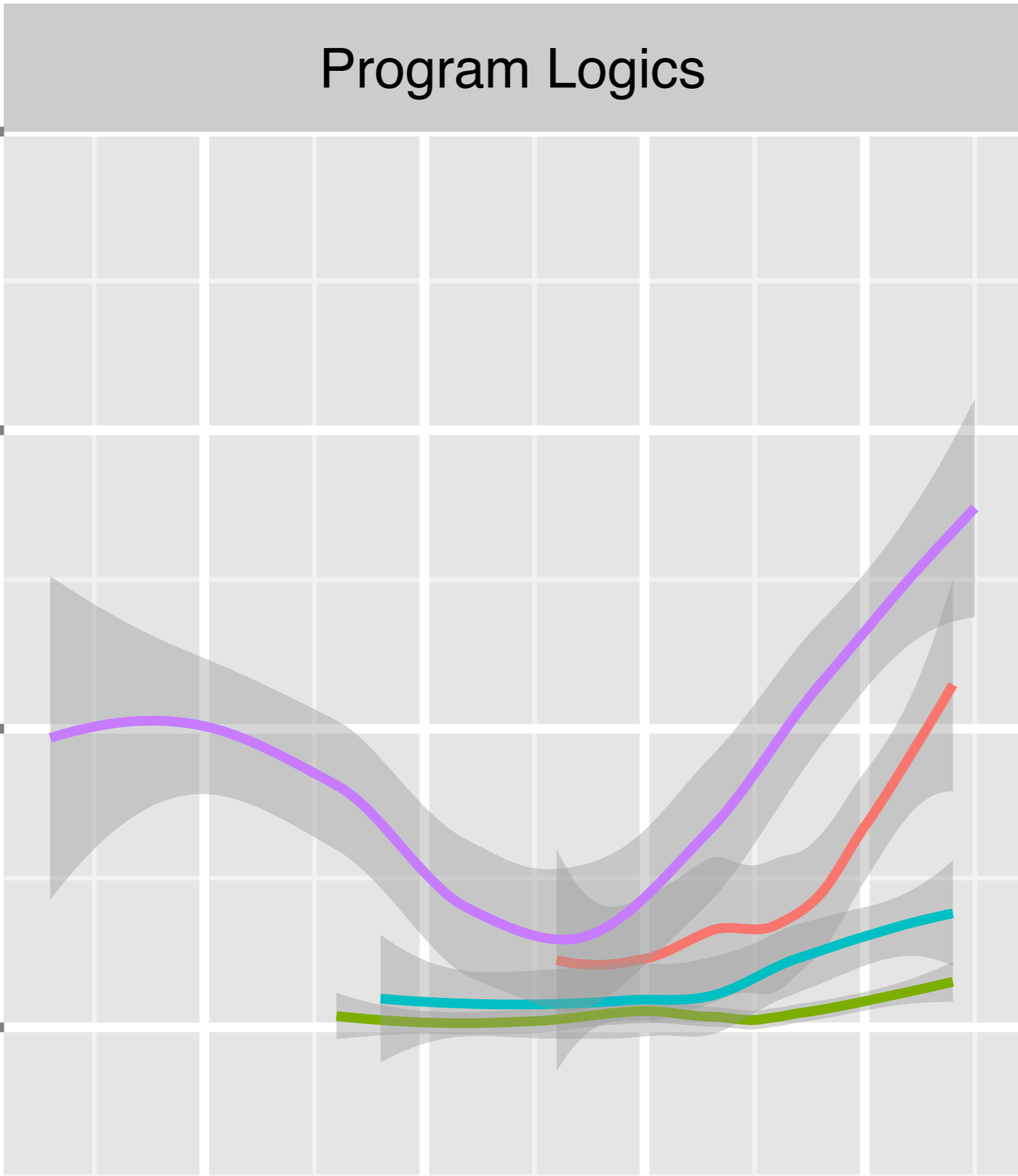
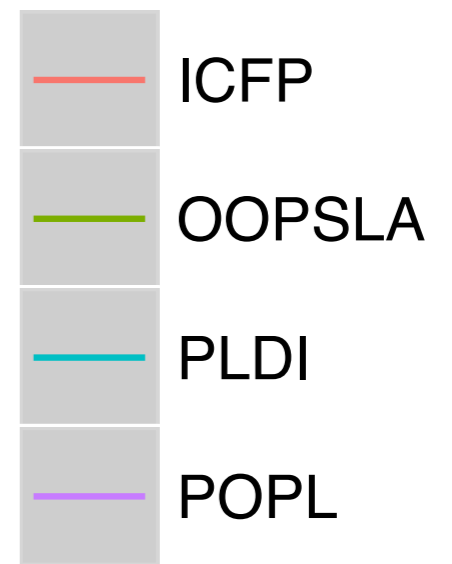
30

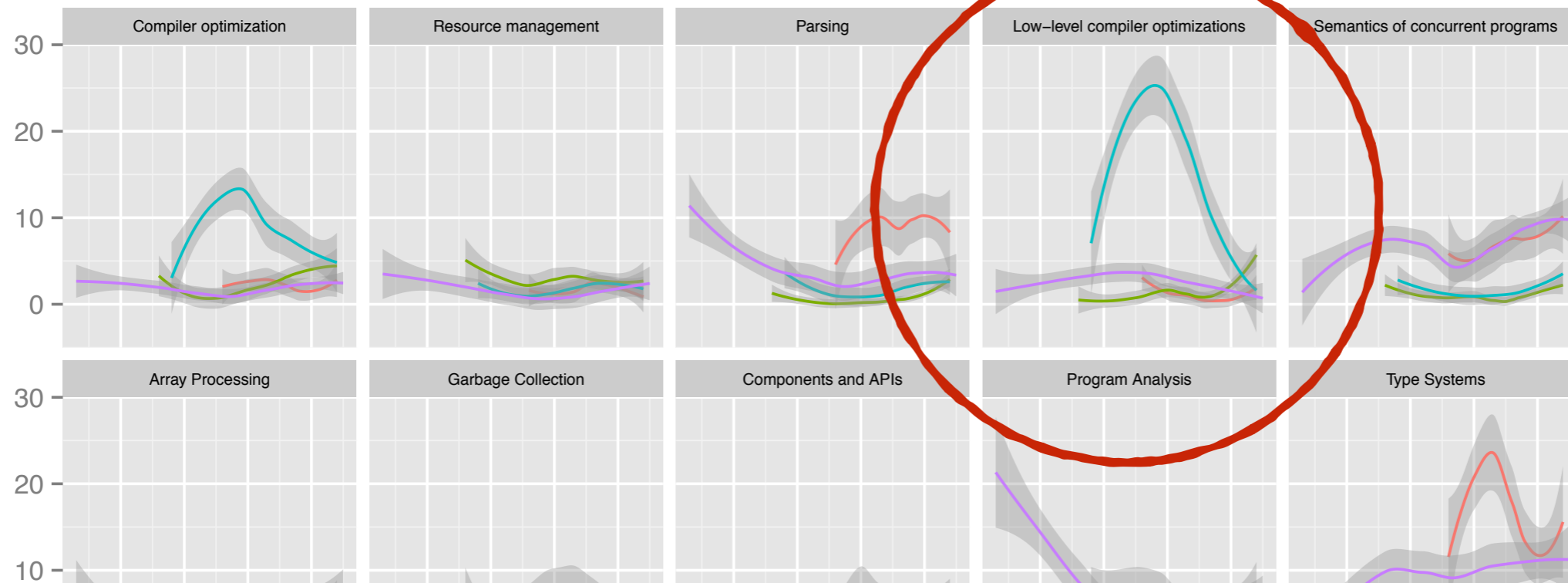
20

10

0

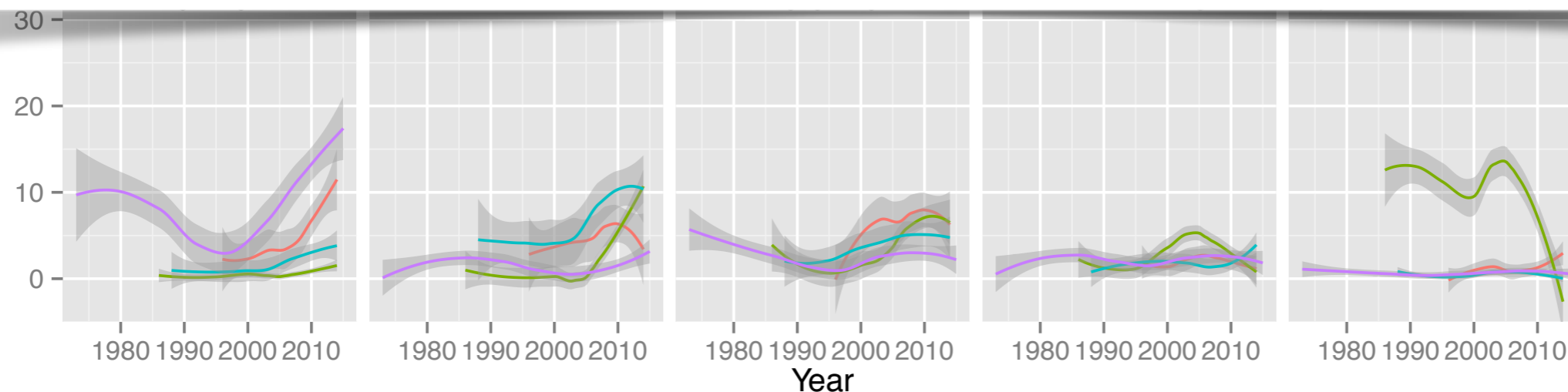
1980 1990 2000 2010



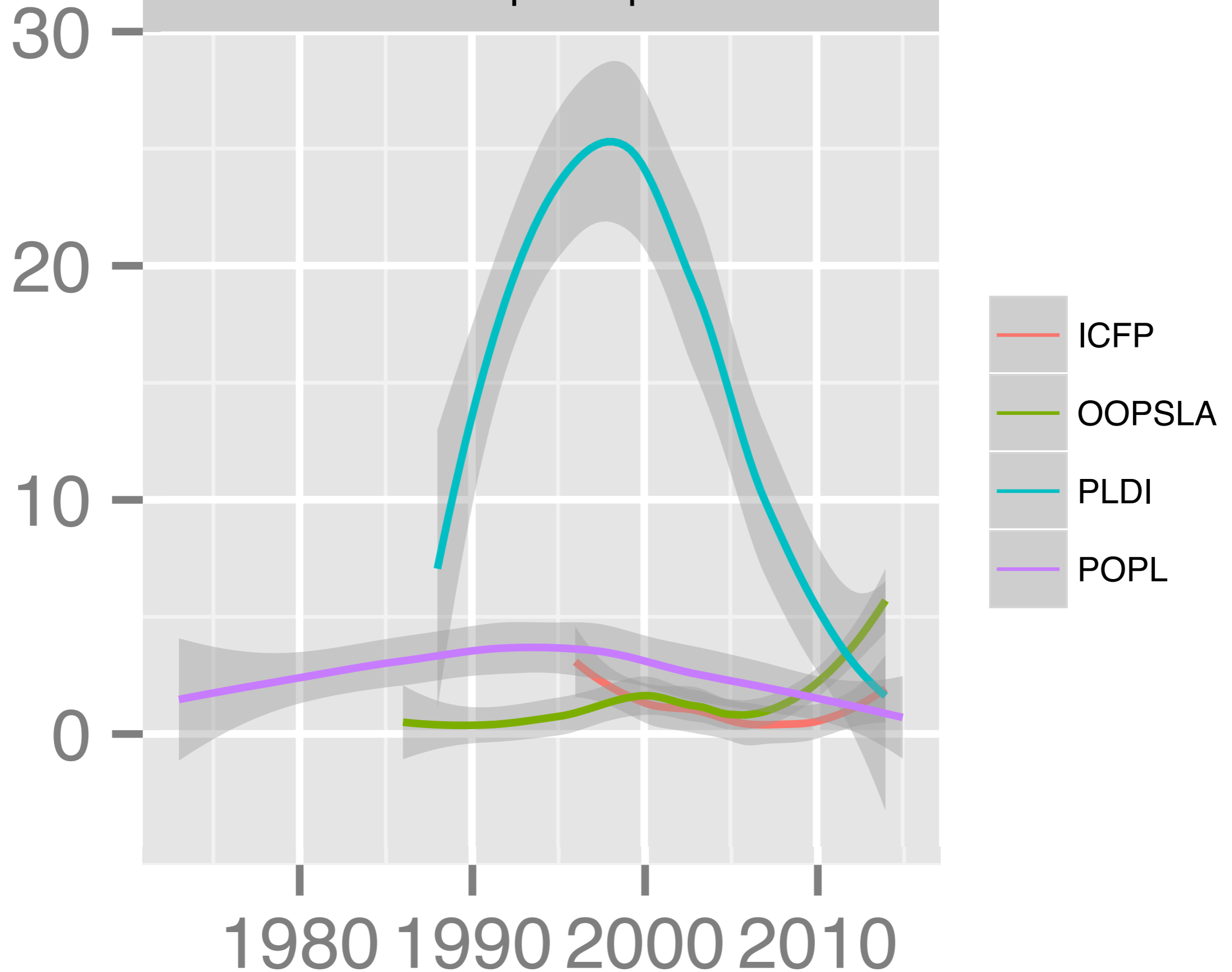


How has **PLDI** changed over time?

Per “Future of PLDI” session in Edinburgh, what is the state of the community?



# Low-level compiler optimizations



# Topic names for k=20, full text

Data-driven  
optimization

Abstract  
interpretation

Object-  
orientation

Code generation

Data-structure  
correctness

Languages and  
control

Security and  
bugfinding

Processes and  
message passing

Garbage  
collection

Parallelization

Program  
transformation

Dynamic analysis

Low-level  
systems

Design

Program analysis

Proofs and  
models

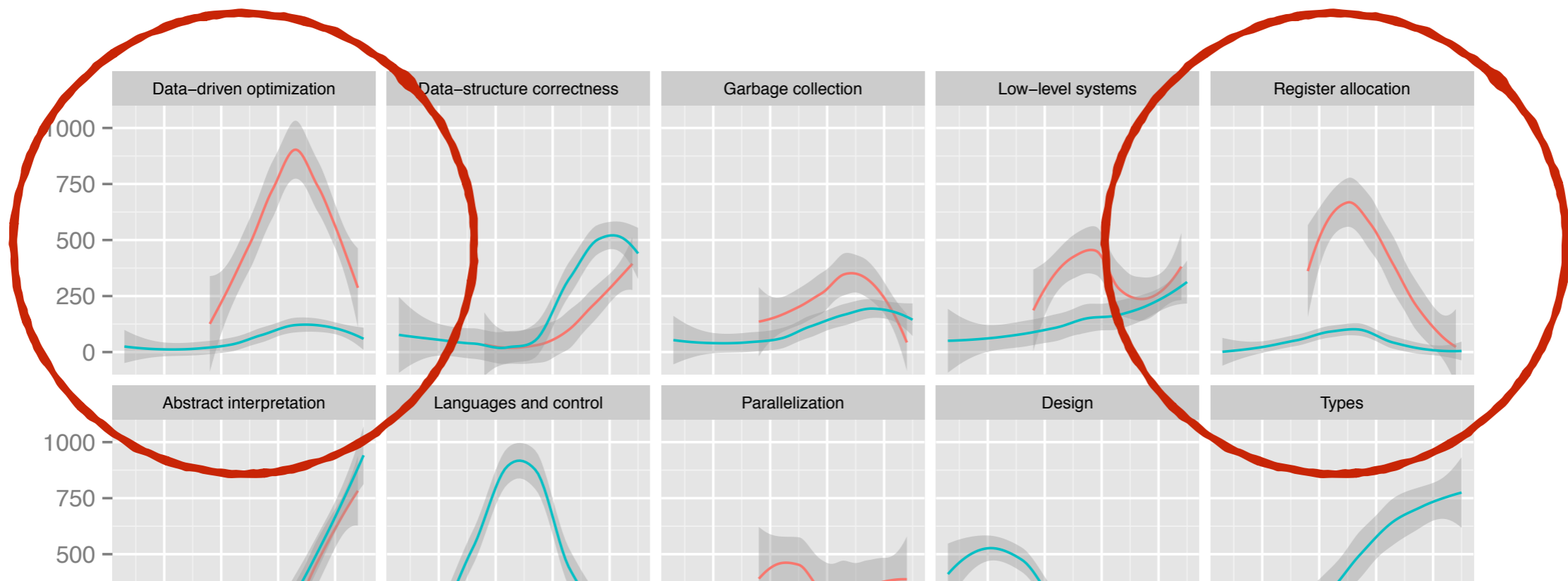
Register  
allocation

Types

Concurrency

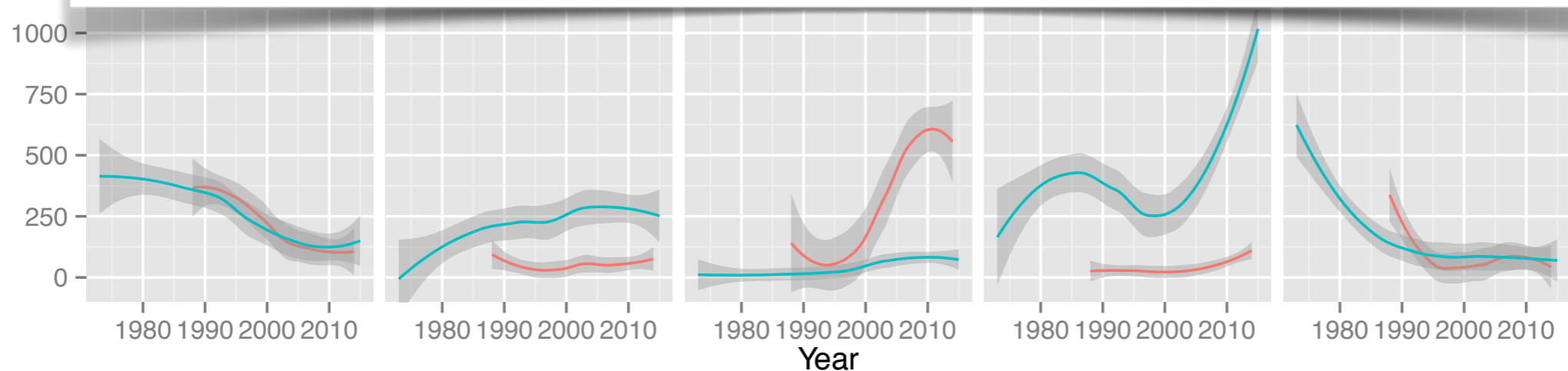
Parsing

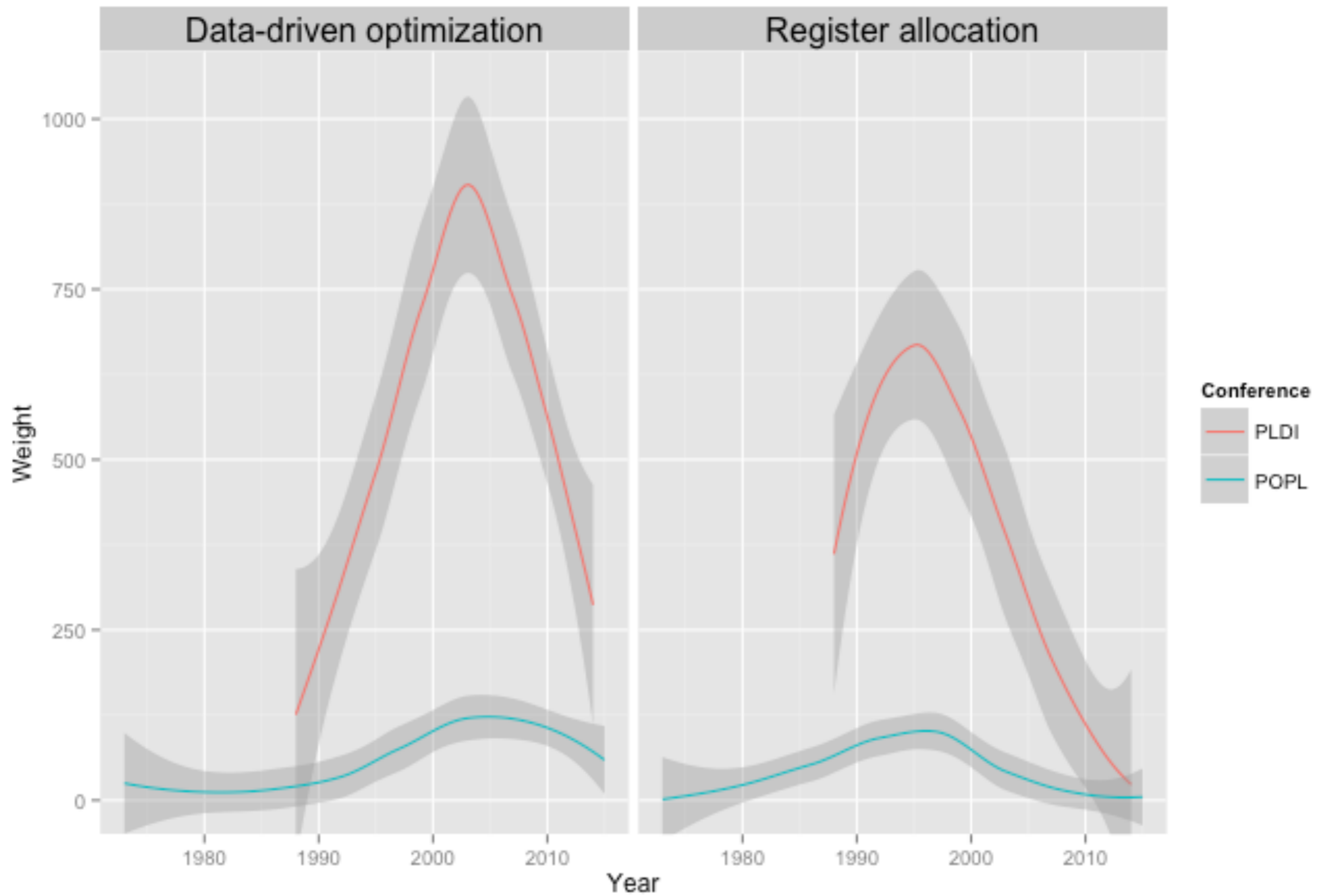


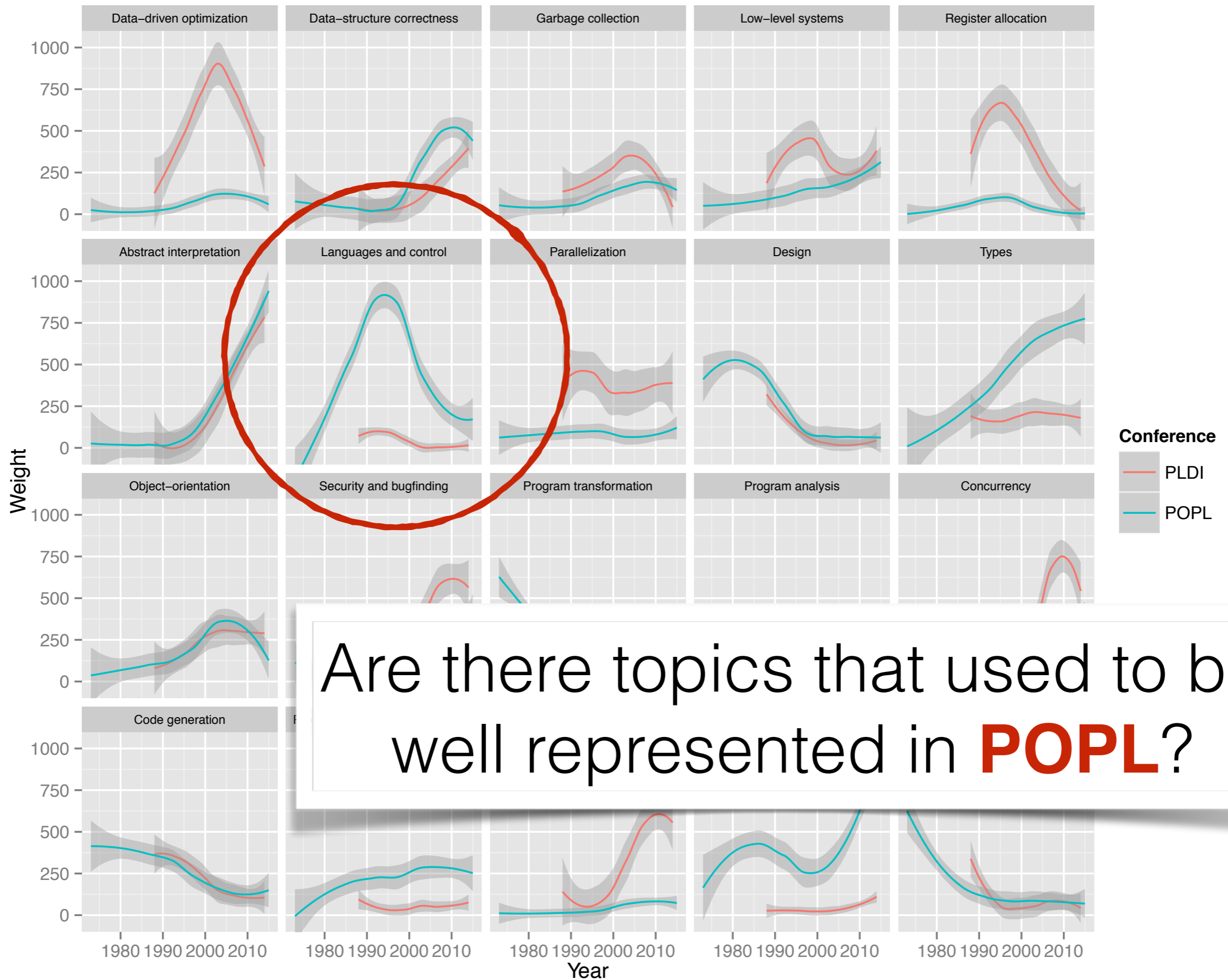


How has **PLDI** changed over time?

Let's compare **PLDI** and **POPL**,  
using our fulltext corpus.

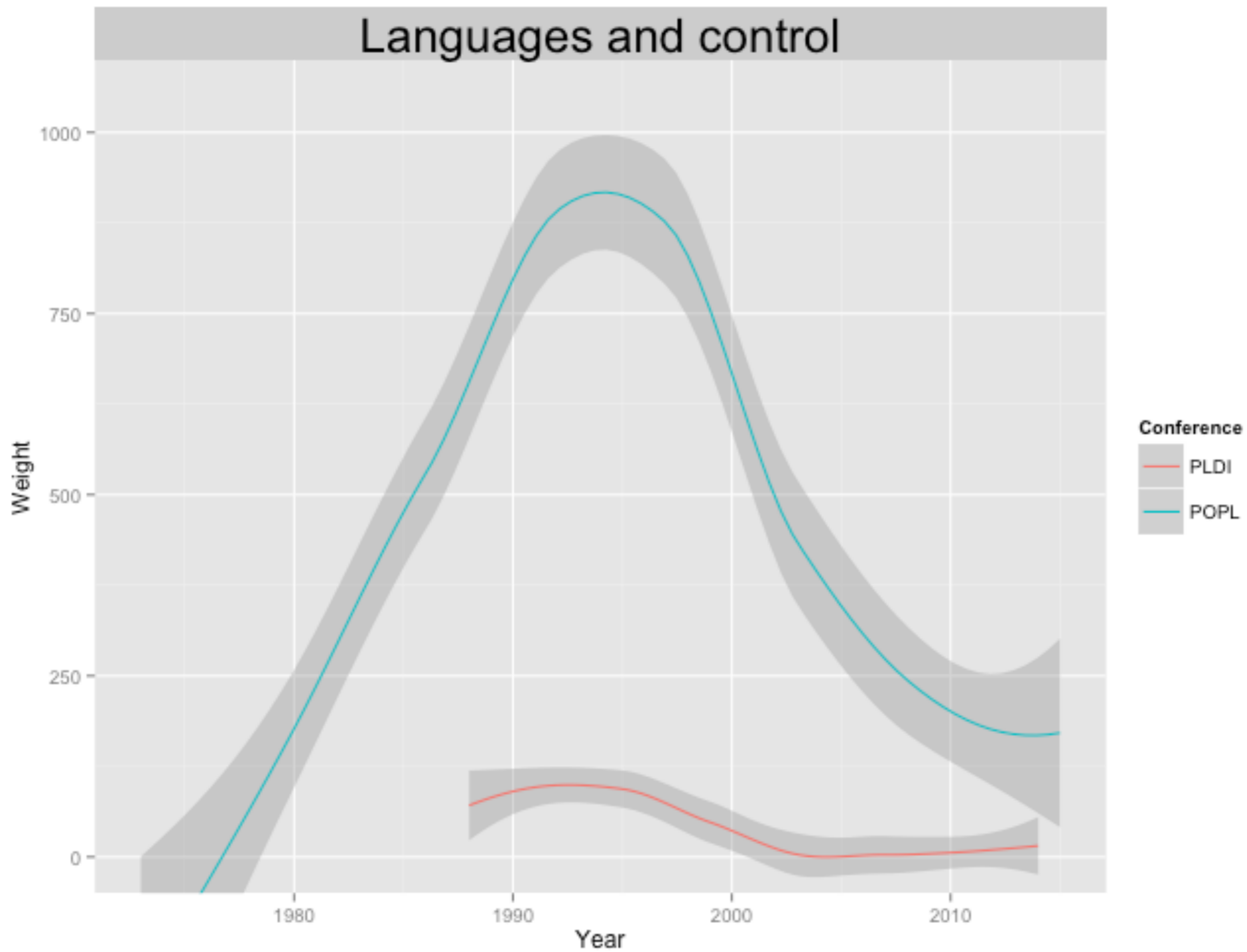




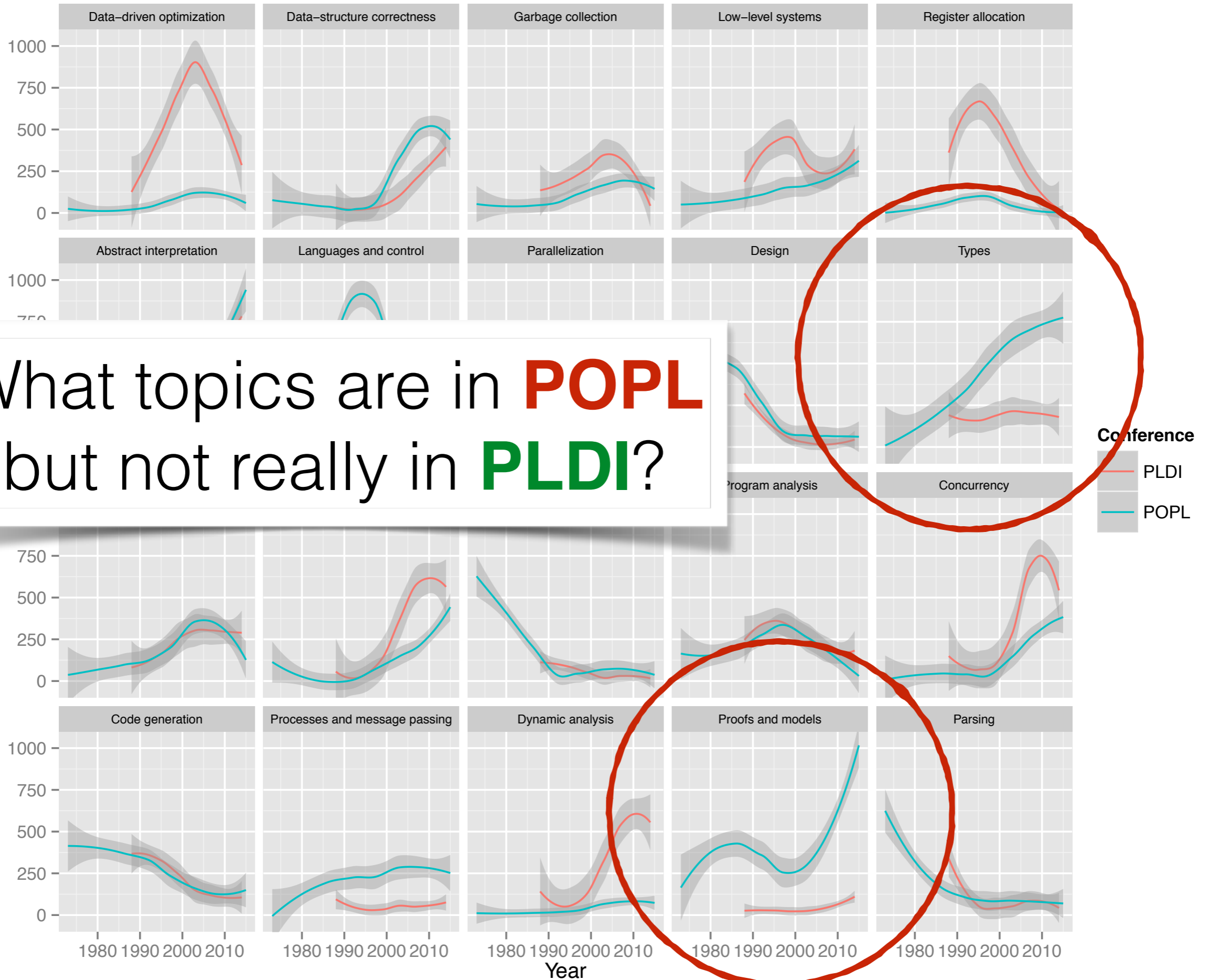


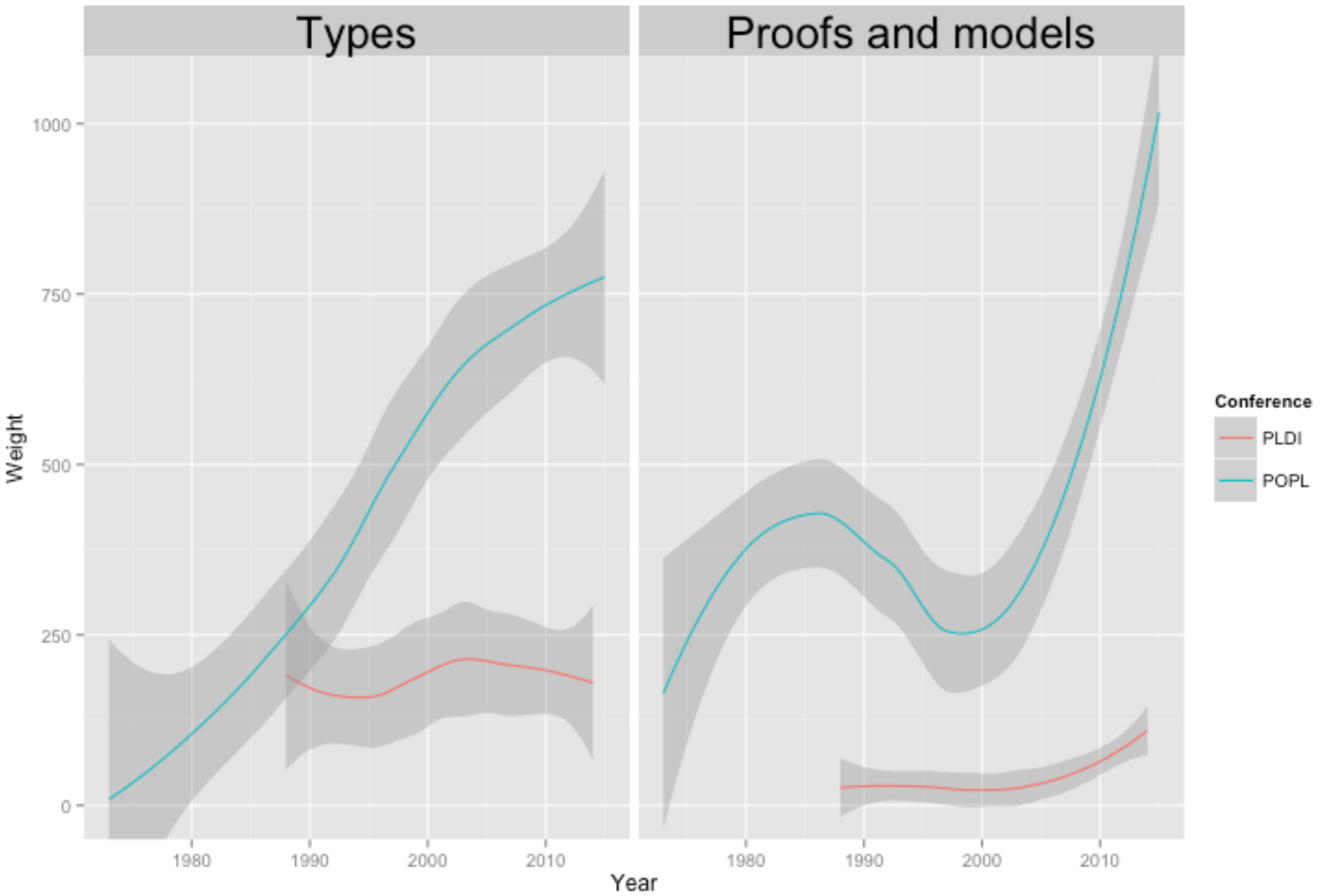
Are there topics that used to be well represented in **POPL**?

# Languages and control

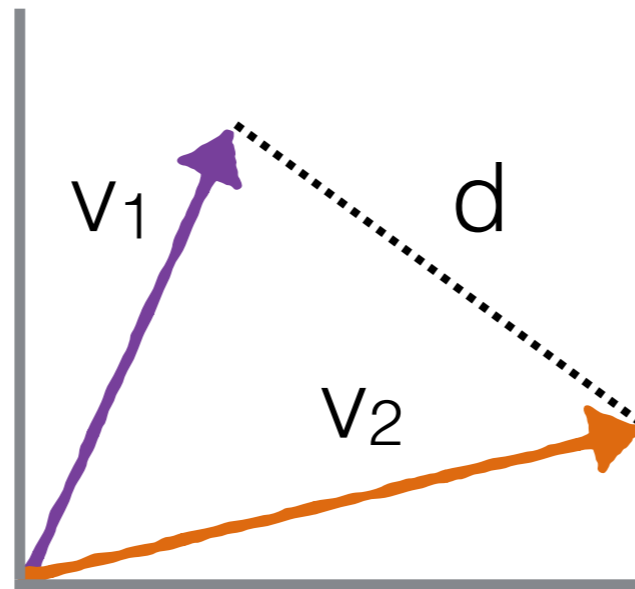


What topics are in **POPL**  
but not really in **PLDI**?



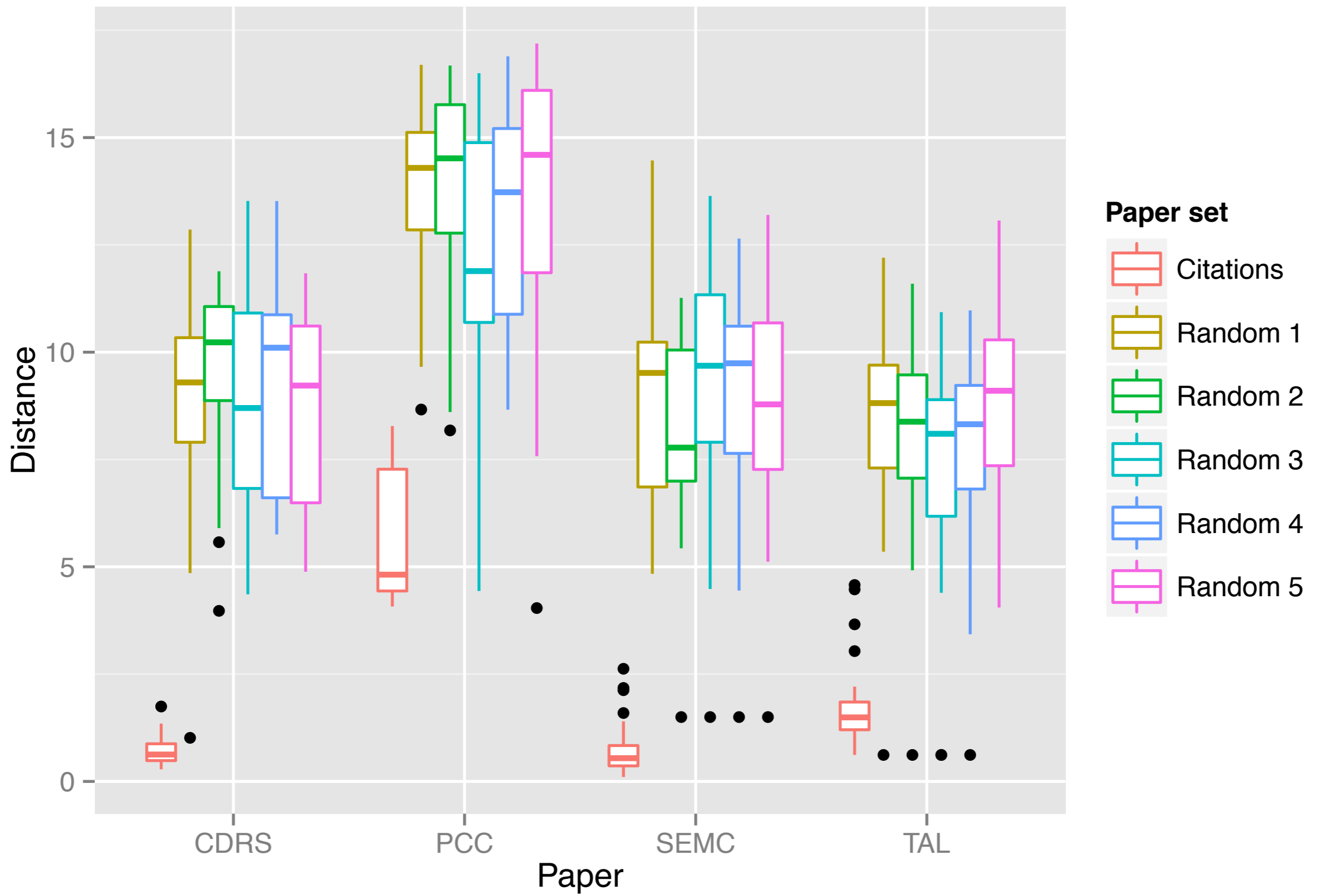


# Comparing documents



Are papers with **close** topic vectors **related**?

Measure distance using *Symmetrized KL divergence*, which gives less weight to dimensions with small magnitude.





<http://tmpl.weaselhat.com>

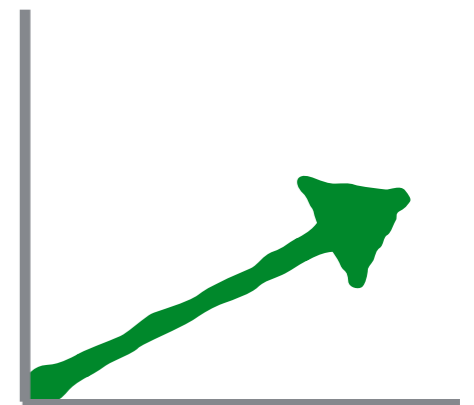
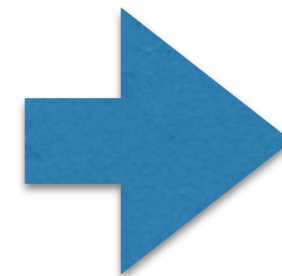
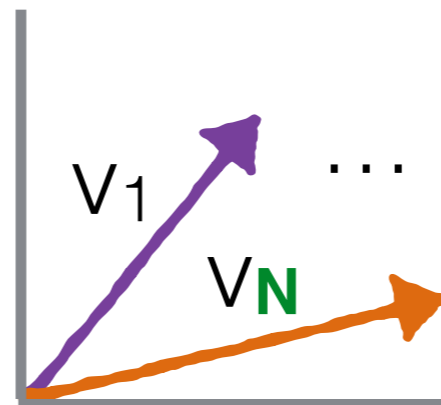


# Ideas and plans

Beginning of a new project

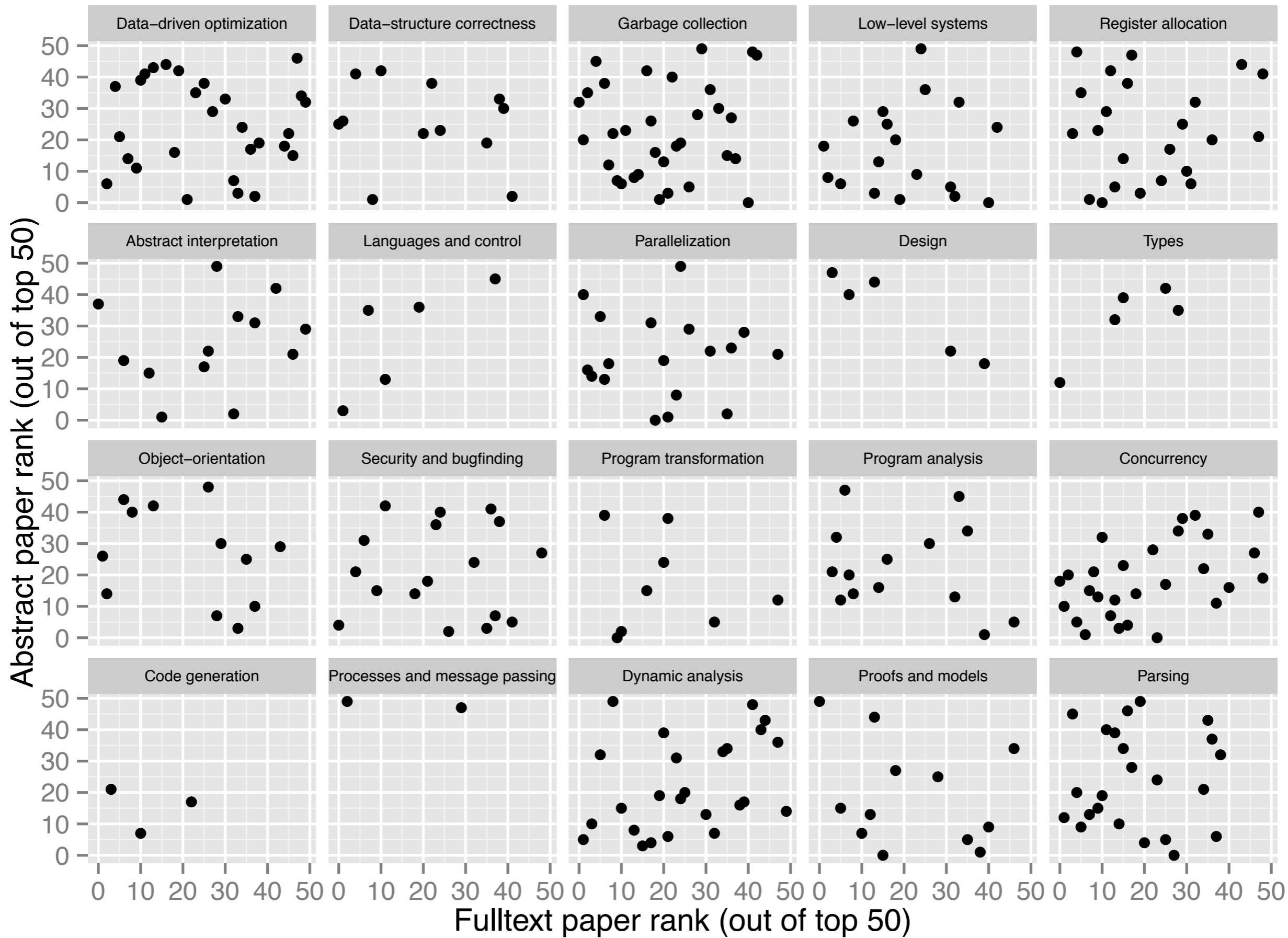
What do *you* think we should do?

Models for researchers



# Limitations/problems

- ACM DL is missing data
  - No programmatic access
- Unclear choices about models
  - Abstracts or fulltext?  $k=20$ ?  $k=30$ ?  $k=200$ ?
  - Which documents should 'seed' LDA?



(More) Questions?