### Time-Related Patterns of Schema Evolution

### EDBT 2025

Panos Vassiliadis, Alexandros Karakasidis http://www.cs.uoi.gr/~pvassil/projects/schemaBiographies/



University of Ioannina, Hellas



University of Macedonia, Hellas

# What are the laws of database schema evolution?



What are the laws of database (schema) evolution?

- How do databases change?
- In particular, how does the schema of a database evolve over time?



- Long term research goals:
  - Are there any "invariant properties" (e.g., patterns of repeating behavior) on the way database (schemata) change?
  - Is there a theory / model to explain them?

Why care for the laws/patterns of schema evolution?

- Scientific curiosity!
- Practical Impact: DB's are dependency magnets. Applications have to conform to the structure of the db...
  - typically, development waits till the "db backbone" is stable and applications are build on top of it
  - slight changes to the structure of a db can cause several (parts of) different applications to crash, causing the need for emergency repairing



### The dreaded schema evolution



- Stonebraker at al., CACM 60, 1 (2017): "In a survey of 20 database administrators (DBAs) at three large companies in the Boston area, we found that . . . , DBAs try very hard not to change the schema when business conditions change, preferring to "make things work" without schema changes. If they must change the schema, they work directly from the relational tables in place. "
- Limoncelli CACM 62, 1 (2019): "When the software is tightly coupled to the database schema it becomes impossible to perform software upgrades that require a database schema change. If you first change the schema, the instances will all die or at least get confused by the change; . . . Why not upgrade the instances first? Sadly, as you upgrade the instances' software one by one, the newly upgraded instances fail to start as they detect the wrong schema. You will end up with downtime until the schema is changed to match the software"

#### FOSS is an opportunity

- Historically, nobody from the research community had access + the right to publish to version histories of database schemata
- Open source tools internally hosting databases have changed this landscape:
  - not only is the code available, but also,
  - public repositories (git, svn, ...) keep the entire history of revisions
- We are now presented with the opportunity to study the version histories of such "open source databases"



From our ICDE 2021 data set collection process that poked 327 repo's....

- 70% of the projects, demonstrated total absence or very small presence of change.
- Out of the 327 repositories that we cloned,
  - 132 (40%) had <u>a single commit</u> for their schema (i.e., no change) whatsoever,
  - 34 (10%) had more than 1 commits, but zero changes at the logical-level schema, and,
  - 65 (20%) were almost frozen (with less than 4 active commits and 10 modified attributes).
- We have called this phenomenon gravitation to rigidity in our past research [IS15, IS17, JoDS17]

http://www.cs.uoi.gr/~pvassil/projects/schemaBiographies/7

Within the 195 that we cloned

#### Gravitation to Rigidity: the reluctance to evolve the schema is omnipresent, stronger than the tendency to evolve, and grows stronger over time!



### Goal of the paper

- During the last 10 years, we have seen several studies of the characteristics of relational schema evolution with a focus on the volume, change-type breakdown and impact to surrounding queries – however, to the best of our knowledge, none has come up with patterns of how schemata evolve over time.
- The goal of this paper is to investigate the timing nature of schema evolution and extract patterns of how schemata evolve over time.

#### Our core contribution

- The major contribution of this paper is that, for the first time in the related literature, a set of time-related patterns on how schema evolution unfolds in time is introduced.
- To achieve this, we had to overcome several limitations:
  - Need to use a large corpus of schema histories to obtain generalizable results
  - Absence of any prior knowledge on how these patterns would look like
  - No established method in our discipline for extracting patterns from schema histories
- To address the volume challenge, we employ a large corpus of 195 relational, logical-level schema histories, out of which we extracted 151 schema histories of length higher than 12 months for further study.
- To address the methodological gap, we have employed methods from empirical software engineering, and iteratively grouped the collected schema histories in patterns of similar time evolution in a qualitative fashion, later to be quantitatively verified, too. We also verified the cohesion and disjointedness of the patterns.

### Time-Related Patterns of Schema Evolution

- We extracted...
- ... 8 Patterns of Change in Time...
- ... organized in 3 Families
- We verified the results of pattern extraction wrt common sense, generalization, disjointness, cohesion, completeness
- We related the patterns to other properties of schema change



### Contributions 1/2

- The core contribution of this paper is the identification of 8
   patterns of schema change, organized in 3 families, on how
   developers and data curators regulate schema evolution over
   time.
- The patterns essentially reflect a model of how change is done via **two important traits**:
  - aversion to change, practically 2/3 of the corpus, and,
  - observable, regular evolution, in several fashions: rare or dense, yet regular, change (amassing to 25% of the corpus), and, surprisingly, an 11% of the corpus with late change too.
- We verified the results of pattern extraction common sense, generalization, disjointness, cohesion, and completeness

### Contributions 2/2

- Other measures of evolution: Although all patterns have similar PUP, the Smoking Funnel and Regularly Curated projects start bigger and contain larger schema evolution activity than the rest of the projects who start small and typically show lower values of change.
- **Change types**: The projects of the change-averted patterns come with small change, frequently being zero, and an inclination towards expansion. The rest of the patterns come with higher volumes of change, and a variety of change types, mostly towards expansion. Both expansion and maintenance are performed with the granule of change being mostly the entire table.
- **Point of Birth**: 34% of the schemata are born in M0, 60% in the first 6 months and 68% in the first 12 months
- **Prediction**: The point of schema birth, gives an early, coarse indication of the subsequent evolution: if born in M0 or after the first year, the schema has a strong inclination towards rigidity (75% and 64% resp.); birth within the first year however, gives a 53% probability, respectively.

Experimental Setup and Nomenclature

### Scope of the study

- We are interested in the monitoring of the evolution of the <u>logical-level</u> <u>relational</u> schema for <u>significant</u> <u>Free Open Source</u> <u>Software</u> projects, hosted in GitHub.
- We are not covering or generalizing to
  - ... proprietary schemata outside the FoSS domain,
  - ... conceptual or physical schemata,
  - ... non-relational schemata, e.g., XML, JSON, ...

### Extraction Process: aimed for massi

#### history collection at very large numbers

- We use the ICDE 2021 Schema Evolution data set (*P. Vassiliadis. Profiles of Schema Evolution in Free Open Source Software Projects. ICDE 2021*).
- Queried the GitHub Activity Data dataset from Google Cloud BigQuery (a 3TB+ dataset that contains a full snapshot + the commits of more than 2.8 million open source GitHub repositories) for repos having .sql files → 133K repo's
- Joined this with Libraries.io dataset (metadata for > 2.7M FOSS prj's) and
- filtered for
  - original repositories, with more than 0 stars, more than 1 contributor
- excluding
  - all files with 'test' or 'demo' or 'example' in the path
  - instances of multiple appearances of a DDL file for >1 vendors
  - multiple DDL's (file-per-table mode), incremental maintenance, vendor X language Cartesian Products

=> 365 candidates, locally cloned, cleaned from empty .git, .sql files with no CREATE TABLE statements, ..., which eventually led to the final data collection.

=> RESULT: 327 histories out of which

132 (40%) with just a single commit ⇔ never changed ANYTHING)!!! 195 histories with at least an extra commit, which we subsequently used



### We work with <u>significant</u> projects

- In whatever follows, remember that we have not selected just any random project, but rather,...
- we intentionally restricted our scope to original, stared projects, where people were actually contributing effort to develop and maintain.
- Overall, 65% of projects spanned more than 24 months and 77% more than a year.

### Post-identification workflow for each of the 195 projects



Eventually, for each project, we ended up with the automatically extracted + time series of changes + collected stats on timing, schema size & activity + extra statistics manually extracted







Project #Active commits #Areeds postV0 #ATurf postV0 Turf Ratio Turf absence / presence **DurationInDays DurationInMonths DurationInYears** #Commits #Tables@Start #Tables@End #Attrs@Start #Attrs@End **TotalTableInsertions TotalTableDeletions TotalAttrInsWithTableIns TotalAttrbDelWithTableDel TotalAttrInjected** TotalAttrEjected **TatalAttrWithTypeUpd TotalAttrInPKUpd TotalExpansion TotalMaintenance** TotalActivity



# Extracting schema and project histories at the month level

- The ICDE21 data set has been grouped by to provide Monthly Schema Activity (sum of changes performed to the schema, with #changed attributes being the unit of measurement)
- The Monthly Project Activity based on the project histories for each projects, was obtained by:
  - Locally cloning projects from Github;
  - For each commit, the names of the changed files, the date, some extra info on the authors and their msg's;
    - git log -name-status -nomerges -date=iso
  - Counted # changed files per commit ;
  - Group by month to produce Monthly Project Activity

### Pattern Extraction Methodology

- 1. We have excluded all projects with a life time less or equal to 12 months: 151 projects.
- 2. We **manually** searched for patterns of the schema line and annotated projects accordingly.
  - This process was iterative, in several rounds and based solely on the aforementioned visual representation of the cumulative progress of schema evolution.
  - Why intentionally manual? Typical in research design s.t. a golden standard of meaningful, humanly-verified groups is attained first, and then checked on the rest of the properties
  - See the paper for pointers on **Grounded Theory** for iteratively extracting patterns out of data
- **3. Quantitatively verified the disjointness ,cohesion and completeness** of the patterns and grouped patterns in larger families.
- 4. Quantitatively analyzed how patterns related to other properties of schema evolution

All data, results, charts and auxiliary analyses are available at :

https://github.com/DAINTINESS-Group/Schema\_Evolution\_Datasets/

### Nomenclature



Schema Update Period: time span between O<sup>th</sup> (originating v.) and last commit for schema updates Project Update Period: resp., for all project updates Schema Expansion: attr's born with new table, injected to existing tables Schema Maintenance: att's deleted with deleted table, ejected from surviving table, data type change, PK change Schema activity = Expansion + Maintenance Unit of measurement: #affected attributes

**Horizontal axis**: time as a percentage of a project's life.

- Vertical axis: cumulative progress as a percentage of the total amount of evolution activity, for
- (a) the schema (dotted, blue line)
- (b) (b) the source code (solid, green line).
- Top-band: 90% of total activity
- **Growth period**: between schema birth and attainment of top-band
- Vault: when the transition between schema birth and top-band takes less than 10% of the total time.

# Quantization of the Measures of Schema Evolution

%Schema					
Activity at	Low	Mic	ldle	High	Full
Birth	≤0.25	(0.25	.0.75]	(0.751)	1
Time of Birth	V0	Ea	rly	Middle	Late
(%PUP)	0	(0 (	0.25]	(0.250.75]	> 75%
Time of					
entering top	V0	Ea	rly	Middle	Late
band (%PUP)	0	(00	0.25]	(0.250.75]	> 75%
Interval		Soon Fair			Very
(%PUP) (birth	Zero	(0	(0.1	Long (0.35	Long
top-band)	0	0.1]	0.35]	0.75]	> 75%
Interval					
(%PUP) (top-	Soon	Fair		Long	Full
band end]	≤0.25	(0.250.75]		(0.751)	1
Active					
months as				Fair (0.2	High >
%growth	Zero 0	Few (0 0.2]		0.75]	75%
Active					
months as				High (0.08	Ultra >
%PUP	Zero 0	Fair (0 0.08]		0.5]	50%

We show (a) how the different metrics have been quantized in labels, and (b) #projects that pertain to each of the produced labels.

- We quantized the different metrics, to be able to perform subsequent analyses
- The quantization was made having in mind the goal to produce general classes of behavior rather than overfit the data set which results in few exceptions in the patterns exactly because the labeling was not done with the patterns in mind.
- The decision of limits was made with the goal of providing labels that (a) are reasonable, and, (b) coarsely divide the metric into labels with similar amounts of projects pertaining to them.
- The extreme values (zero and one) have particular semantics in most cases, too.

### Statistical Study

Overall: statistical characteristics of the introduced measures, before studying the corpus breakdown in patterns

### Volume of Birth (%Total Change).

Volume of	Low	M	iddle	High (0.75, 1)	Full	
Birth (%Total Change)	<=0.25	(0.2.	0.75]	(0.751)	1	
chunge)	(16)		(52)	(44)	(39)	
Time Point of	$V_p^0$	Early		Middle	Late	
Birth (%PUP)	0	(0 .	. 0.25]	(0.250.75]	>75%	
	(52)		(53)	(33)	(13)	
Time point of	$V_{p}^{0}$	E	arly	Middle	Late	
reaching Top Band (%PUP)	0	(0 .	(0 0.25]		>75%	
	(23)	(41)		(47)	(40)	
Interval	Zero	Soon	Fair	Long	Very Long	
(%PUP) (birth top-band)	0	(0 0.1]	(0.1 0.35]	(0.35 0.75]	>75%	
1 /	(62)	(26)	(27)	(23)	(13)	
Interval	Soon	]	Fair	Long	Full	
(%PUP) (top- band end]	<=0.25	(0 0.25]		(0.751)	1	
build elitaj	(40)	(48)		(40)	(23)	
Active months	Zero	Few		Fair	High >75%	
as %growth	0	(0 0.2]		(0.2 0.75]	1 HgH ~7.5%	
5	(98)	(22)		(22)	(9)	
Active months	Zero	]	Fair	High	Ultra >50%	
as %PUP	0	(0 0.08]		(0.08 0.5]	21114 - 0070	
	(98)	(20)		(33)	-	

Reading from right to left, we see that out of 151 projects,

 39 projects reach Full (100%) activity at schema birth, and

83 projects overall reach
 High or Full activity at
 schema birth.

 Overall, more than half of the projects exceed 75% of total activity at schema birth.

### Timepoint of Birth (%PUP)

Valence of	Low	Middle			High	Full		
Birth (%Total	<=0.25	(0.250.75]		(0.751)	1	•		
Change)		()		(	()			
	(16)	(52)		(44)	(39)			
Time Point of	$V_p^0$	Early (0 0.25] (53)		Middle	Late			
Birth (ZPLIP)	Ō			(0.250.75]	>75%			
	(52)			(33)	(13)			
Time point of	$V^0_p$	E	arly		Middle	Late		
Time point of	Ô	(0.	. 0.25]		(0.250.75]	>75%		
Band (%PUP)							•	
	(23)	(41)		(47)	(40)			
Intonnal	Zero	Soon	Fai	r	Long	Very Long		
(or DLID) (birth	0	(0 0.1]	(0.1 (	0.35]	(0.35 0.75]	>75%		
(%FUF) (bittil							•	
top-band)	((n))	(26)	(07	n)	(02)	(12)		
	(62)	(20)	(27	)	(23)	(13)		
Interval	Soon			Long	Full			
(%PUP) (top-	<=0.25	(0 0.25]		(0.751)	1			
band end]								
	(40)	(48)		(40)	(23)	1		
Active months	Zero	Few (0 0.2]		Fair	High >7507			
Active months	0			(0.2 0.75]	1 ligii >7 5%	1		
as //growin	(09)				(00)	(0)		
	(98)	(22)		(22)	(9)			
Active months	Zero	Fair		High	Ultra >50%	0		
as %PUP	0	(0.	. 0.08]		(0.08 0.5]			
	(98)	(20)		(33)	_	0		

- A large majority of schemata is born early: 52 schemata (one third of the population) are born in V<sup>0</sup>.
- Two thirds of the projects (105 projects) see schema birth at V<sup>0</sup> or before 25% of the PUP.
- 74 schemata (half the corpus) are born in the first 10% of time.
- As the distribution of points of birth follows a power-law shape, the rest of the time points form a long tail.



### Time point of Entering the Top Band (%PUP)

Volumo of	Low	M	iddle	High	Full	]
Volume of Dirth (gratal	<=0.25	(0.25	50.75]	(0.751)	1	•
O(1 - 1)						
Change)	(1.1)		(= 0)	(11)	(2.2)	
	(16)		(52)	(44)	(39)	
Time Point of	$V_p^0$	Early Middle		Late		
Birth (70110)	Ô	(0.	. 0.25]	(0.250.75]	>75%	
Biitii (%F UF)	(50)		(50)	(0.0)	(10)	
r	(52)		(53)	(33)	(13)	
Time point of	$V_p^0$	E	arly	Middle	Late	
reaching Ton	0	(0.	. 0.25]	(0.250.75]	>75%	
Band (7 DUD)						
Dallu (%r Ur)	(0.0)		(41)	(477)	(40)	
	(23)		(41)	(47)	(40)	
Interval	Zero	Soon	Fair	Long	Very Long	•
(%PUP) (birth	0	(0 0.1]	(0.1 0.35]	(0.35 0.75]	>75%	
ton-band)						
top build)	(62)	(26)	(27)	(23)	(13)	
	(02) Soon	(20)	(27) Foir	(23) Long	(13) Eull	-
Interval	50011	(0  0.25]		(0.75, 1)	1 run	
(%PUP) (top-	<=0.25	(0.	. 0.25]	(0.751)	1	
band end]						
1	(40)		(48)		(23)	
A	Zero	Few		Fair	TT: 1 ===	
Active months	0	(00.2]		(0.2 0.75]	High >75%	
as %growth	(2.2)		(0.0)	(2.2)	(2)	
	(98)	(22)		(22)	(9)	
Active months	Zero	Fair		High	Ultra >50%	
as %PUP	0	(0 .	. 0.08]	(0.08 0.5]	5111a - 5070	
	(98)	(20)		(33)	-	

64 projects (42%) reached the top band immediately at V<sup>0</sup> (23 of them) or before 25% of the PUP. Another 40 projects (26%) came with a late top-attainment at higher than 75% of the PUP. The rest of the projects were spread in the middle half of the PUP.

In other words, projects mostly reach topband soon, although mid-life and late completion of the evolution do exist too.



### Interval Schema Birth-To-TopBand (%PUP)

Volume of	Low	M	iddle	High	Full		Vaulte
Birth (%Total	<=0.25	(0.2	50.75]	(0.751)	1	•	vauits.
Change)							vault, i
	(16)		(52)	(44)	(39)	_	
Time Point of	$V_p^0$	Early		Middle	Late		top-bai
Birth (%PUP)	0	(0 .	0.25]	(0.250.75]	>/5%		with 62
	(52)		(53)	(33)	(13)		
Time point of	$V_p^0$	E	Carly	Middle	Late		115 /70
reaching Top	0	(0 .	0.25]	(0.250.75]	>75%	•	112 (/:
Band (%PUP)	()		(				less tha
r	(23)	C	(41)	(47)	(40)	_	
Interval	Zero 0	(0 0.1]	$(0.1 \dots 0.35]$	(0.35 0.75]	>75%	•	The ris
(%PUP) (birth		(0 0.1]	(0.1 0.00]	(0.00 0 0]			
top-band)	(62)	(26)	(27)	(23)	(13)		therefo
12			BirthToTo	opBand interval	(%PUP)	1	
1.12							
1							
0.8							
0.6							
0.0					_		
0.4					-		
0.2						eeeeeeeeee	and the second sec
V12					A COCCORDE COLOR		

**Vaults**. 88 / 151 projects (58%) had a single vault, i.e., an interval from schema birth to top-band that was less than 10% of the PUP, with 62 of them in zero time.

- 115 (75%) of the projects had an interval of less than 35% PUP.
- The rise from schema birth to the top is therefore mostly fast.

4 7 10 13 16 19 22 25 28 31 34 37 40 43 46 49 52 55 58 61 64 67 70 73 76 79 82 85 88 91 94 97 100103 106 109 112 115 118 121 124 127 130 133 136 139 142 145 148 151

#### Vaults and Active Growth Months



- Active Growth Months. 98 schemata (2/3 of the population) had zero active months in the growth period, from schema birth to the top.
- Another 17 had exactly one month active. Combined, 115 projects, 76% of the population had no more than 1 month of activity from schema birth to top.
  - If there should be a single metric of
    evidence for the aversion to change,
    the super-focused nature of schema
    evolution, and, schema stabilization,
    Active Growth Months would be the
    one.

### Correlations



Aversion to change is predominant: the longer it takes to reach the top-band, the more inactive months you include, i.e., people prefer clustered groups of schema changes rather than constant incremental maintenance. The different measures of Schema Evolution are correlated to some extent.

- The ActiveGrowthMonths (used by default in our analysis) is very tightly related to its normalized versions as percentages of the Project Update Period and the Growth Period.
- Volume of Change at birth and the Interval from Birth To Top as %PUP are strongly related to the months without change in the growth period
- Point of Top-Band Attainment and the Interval from the Top-Band to the End of the project are extremely strongly anticorrelated and strongly correlated to Point of Birth (all of this is totally expected)

### Summary of simple stats

- Birth is mostly done early
  - Two thirds of the projects (105 projects) see schema birth at V0 or before 25% of the PUP.
  - 74 schemata (half the corpus) are born in the first 10% of time.
- Change is mostly sharp in time:
  - 115 projects, 76% of the population had no more than 1 month of activity from schema birth to top.
  - 88 / 151 projects (58%) had a single vault, i.e., an interval from schema birth to top-band that was less than 10% of the PUP, with 62 of them in zero time.



### The Patterns

8 Patterns of Change in Time... ... organized in 3 Families



Smoking

0.8 0.9

**Funnel** 

The "Be Quick or Be Dead" family of patterns constitutes a family of very focused change very close to the point of schema birth - the only difference of the involved patterns is when schema birth takes place.





The "Stairway to Heaven" family of patterns: both patterns, involve a fairly regular pattern of change, with change steps distributed across time.

Although **different in the change rate**, both patterns refer to projects that do not reach the top band in a single shot, but, **progressively climb to the top-band over a long period of time**.



The **"Scared to Fall Asleep Again**" family of patterns: the two patterns, although very different in their characteristics, resemble in that they **include projects where the change is not focused in a single point**, and **happens towards the end of the lifetime of the project**.


### Flatliners

Everything happens at birth



#Prjs	Born	WhenHow long isreaches topgrowth?band?(middle life)		How long a tail?	# Active months at growth	
Out of 151	Early <=25%, middle (25% 75%], late > 75%	%,         Early <=25%,         When /how from birth to top Band         How           %75%],         middle (25%75%],         (>90% totAct)?         to en           late > 75%         0, soon<=10%, fair<=35%, long		How long from reach-of-topBand to end? Soon <=25%, fair (25% 75%], long (75% 100%) Full 100%	Growth: [birth-topBand)	
23	V0	V0	Zero	Full	0	
	Born really early	via a single vault	that does 100% of the job			

Exceptions:

## Radical Sign

Born early, very soon freezes totally



#Prjs	rjs Born When How reaches top gro band? (mi		How long is growth? (middle life)	How long a tail?	# Active months at growth	
Out of 151	of 151 Early <=25%, Early <=25%, middle (25% 75%], middle (25% 75%], late > 75% late > 75%		When /how from birth to top Band (>90% totAct)? 0, soon<=10%, fair<=35%, long <=75% very long > 75% time	How long from reach-of-topBand to end? Soon <=25%, fair (25% 75%], long (75% 100%) Full 100%	Growth: [birth-topBand)	
41	VO, Early	Early	zero, soon, fair	Long	0 - 2	
	Born early	mostly via a single vault;	Trip to top is fairly short;	All ends soon, i.e., a long tail		

Exceptions:

## Sigmoid

Born in the middle, sharp vault to top-band Resembles a "pure" sigmoid function, better than any other pattern

(almost all patterns are Exceptions of a sigmoid)



#Prjs	Born	When reaches top band?	How long is growth? (middle life)	How long a tail?	# Active months at growth
Out of 151	Early <=25%, middle (25% 75%], late > 75%	Early <=25%, middle (25% 75%], late > 75%	When /how from birth to top Band (>90% totAct)? 0, soon<=10%, fair<=35%, long <=75% very long > 75% time	How long from reach-of-topBand to end? Soon <=25%, fair (25% 75%], long (75% 100%) Full 100%	Growth: [birth-topBand)
19	Middle	Middle Single vault in the middle	Zero, Soon (~immediately, with a single vault)	Fair	0 - 1

### Late Riser

Born late, sharp vault to top-band



#Prjs	Born	When reaches top band?	How long is growth? (middle life)	How long a tail?	# Active months at growth
Out of 151	Early <=25%, middle (25% 75%], late > 75%	Early <=25%, middle (25% 75%], late > 75%	When /how from birth to top Band (>90% totAct)? 0, soon<=10%, fair<=35%, long <=75% very long > 75% time	How long from reach-of-topBand to end? Soon <=25%, fair (25% 75%], long (75% 100%) Full 100%	Growth: [birth-topBand)
14	Late	Late	Zero, soon	Short	0
			No early or middle life, late single vault		
Exceptions:		(1 exc.: middle)			(1 exc.: 5 months)

### Quantum Steps

A sequence of a few focused steps in the middle...
Two variants, both with few updates:
(a) early start, middle top, and,
(b) middle start, late top



#Prjs	Born	When reaches top band?	How long is growth? (middle life)	How long a tail?	# Active months at growth
Out of 151	Early <=25%, middle (25% 75%], late > 75%	Early <=25%, middle (25% 75%], late > 75%	When /how from birth to top Band (>90% totAct)? 0, soon<=10%, fair<=35%, long <=75% very long > 75% time	How long from reach-of-topBand to end? Soon <=25%, fair (25% 75%], long (75% 100%) Full 100%	Growth: [birth-topBand)
17	V0 or Early	Middle	Fair or Long	Fair	0 – 3
6	Middle	Late	Fair or Long	Soon	0 – 3
23	-				
Exceptions:	-	2 exc.			

## Regularly Curated

Grows "regularly" over time with activity. Reaches top, with activity, in two variants:

- If born early, reaches top middle or late;
- If born middle, reaches top late

Schema line close to the green line of prj evo, occ. with a small tail



#Prjs	Born	When reaches top band?	How long is growth? (middle life)	How long a tail?	# Active months at growth	
Out of 151	of 151 Early <=25%, Early <=25%, middle (25% 75%], middle (25% 75%], late > 75% late > 75%		When /how from birth to top Band (>90% totAct)? 0, soon<=10%, fair<=35%, long <=75% very long > 75% time	How long from reach-of-topBand to end? Soon <=25%, fair (25% 75%], long (75% 100%) Full 100%	Growth: [birth-topBand)	
11	V0 or Early	Middle or Late	Long or very long	Soon	> 3	
3	Middle	Late	Fair or long	Soon	> 3	

### **Smoking Funnel**

Somewhat late birth, with something like a vault

(but not full or super high), and once born,

alive with regular schema updates



#Prjs	Born	When reaches top band?	How long is growth? (middle life)	How long a tail?	# Active months at growth
Out of 151	Early <=25%, middle (25% 75%], late > 75%	Early <=25%, middle (25% 75%], late > 75%	When /how from birth to top Band (>90% totAct)? 0, soon<=10%, fair<=35%, long <=75% very long > 75% time	How long from reach-of-topBand to end? Soon <=25%, fair (25% 75%], long (75% 100%) Full 100%	Growth: [birth-topBand)
7	Middle	Middle	Fair	Fair	> 3
	Born in the middle	and reaches top in the middle	but with some action in the way		thus the activity

Exceptions:

### Siesta

Born early, at a moderate high level,

then goes to sleep for some (significant) time,

then wakes up again



#Prjs	Born	When reaches top band?	How long is growth? (middle life)	How long a tail?	# Active months at growth
Out of 151	Early <=25%, middle (25% 75%], late > 75%	Early <=25%, middle (25% 75%], late > 75%	When /how from birth to top Band (>90% totAct)? 0, soon<=10%, fair<=35%, long <=75% very long > 75% time	How long from reach-of-topBand to end? Soon <=25%, fair (25% 75%], long (75% 100%) Full 100%	Growth: [birth-topBand)
10	V0 or Early	Late	Very Long	Soon	0 – 3
	Early born		with a long sleep in the middle	and some action in the end	i.e., without much ado
Exceptions:			(1 exc.: long)		(2 exc.)

### Traits of schema evolution

- The patterns reflect essentially how change is done, rather than just being statistically-backed project clusters. We observe two important traits:
- The first trait is the aversion to change, aka progressive gravitation to rigidity, meaning that curators avoid change as much as they can, and more often that not, the schema freezes after a few changes. This is a majoritarian trait, and concerns the Be quick or Be Dead family, which involves practically 2/3 of the corpus.
- The second trait concerns a minority of projects whose curation team regularly synchronizes the schema to the surrounding changes, with observable regular schema evolution, coming in several fashions: rare but regular change; densely regular change; and, surprisingly, late change too.
- Overall:
  - the anecdotal evidence of "freeze the schema first; then build all the applications on top of it", although certainly majoritarian as a practice, is only partially corroborated, with the existence of projects that are maintained "regularly" in various fashions
  - The assertion of several works in the related work, that change is frequent, is also mostly disproved – the reason is that, out of necessity, if change is to be studied, it has to be studied

### Pattern Validity

Are these patterns...

- VQ1: genuine and reasonable?
- VQ2: (a) internally cohesive and (b) pairwise disjoint?
- VQ3: generalizable?
- VQ4: complete?

# Validation Questions for this Study

- Genuine?
- Cohesive & disjoint?
- Generalizable?
- Complete?
- VQ1: Are these patterns **genuine** and **reasonable**? How can we guarantee that the separation is not artificial and a-posteriori fitted to the numbers?
- VQ2: Can we claim that the classification of projects into different patterns is producing patterns that are (a) internally cohesive and (b) pairwise disjoint?
- VQ3: How **generalizable**, i.e., how representative of the general behavior of projects, are the results?
- VQ4: Is the taxonomy produced **complete**? How possible is it that other behaviors do exist too?

#### Is what we see ... genuine?

- Genuine?
- Cohesive & disjoint?
- Generalizable?
- Complete?
- Genuine character comes from the process itself: all the statistics-related tasks, like the quantization of the time-related attributes, the pairwise comparison of patterns or the inspection of other attributes were performed only after the bottom-up manual inspection and classification
- The **pattern definitions came also after the manual classification was performed**, and had a very small effect to internal restructurings, just for the Regularly Curated family.
  - In line with **Grounded Theory**, a method suited for situations where patterns have to be extracted from collected data via a manual, iterative generation of patterns via comparing every new data item (here: schema history) to all previous patterns, and adjusting the patterns accordingly (see the paper for references).
- The few **exceptions** that exist are also evidence of the genuine human-based, iterative process.

# Is what we see ... in line with related work?

- "Be quick or be dead" with a single/very few "shot(s)" & small activity volume:
  - too many projects
  - internal separation: only wrt when the schema was born.
- "Stairway to heaven" with progressive evolution at least until midlife,
  - non-trivial fraction of the projects
  - internal separation: wrt the rate of change time-points
- "Scared to fall asleep again" with signs of late change:
  - a small part of the corpus
  - the two patterns differ at the existence of midlife change hiatus

- Studies of the more distant past with a handful of schemata had to pick schemata that have some change (otherwise there is not much to report), so no aversion to change.
- Studies that work with large numbers of schemata, or studies with necessity sampling of projects, clearly show the aversion to change. This means that
  - both Be quick or be dead && Stairway to heaven are in-line with previous research
  - Scared to Fall asleep again is completely new
- The study of nosql schemata: time behaviors that are visually quite close to the patterns reported here (universality of the patterns?)

#### Is what we see ... reasonable?

- Observed desirable properties in our set of patterns
- First, few exceptions do exist. The result of the pattern extraction process was neither perfectly fit to the data, nor with colossal deviations.
- Second, the patterns
  - are disjoint and
  - cover significantly (not fully) the space of possible behaviors.
- Moreover, each pattern came with a reasonable size, neither overly dominating the rest, or being insignificant.



• Genuine?

• Cohesive & disjoint?

Complete?

Generalizable?

#### Disjointness

#### Formal Disjointness: the formal definitions cover disjoint areas in the space produced by the Cartesian Product of values for the defining attributes

		Birth			From Gro	wth To the Top	o Band			Tail
PATTERN FAMILY	PATTERN	Birth Vol Class	Birth Timing Class	Point Top Band Class	Has Single Vault	Interval Birth To TopB Class	Growth Months With Change	Active Pct Growth Class	Active Pct PUP Class	Interval TopB To End Class
Be Quick Or Be Dead	FlatLiner (23)	high, full	V <sub>p</sub> 0	V <sub>p</sub> 0	TRUE	Zero	0	Zero	zero	full
bedd	Radical Sign (41)	low, fair, high, full	$V_p0$ , early	early	TRUE, FALSE	Zero soon, fair	0 – 2 (1 exc)	any	zero, fair	long (1 exc)
	Sigmoid (19)	full, high, fair	middle (2 exc)	middle	TRUE	Zero, soon	0-1	Zero <i>(3 exc)</i>	zero, fair	fair
	Late Riser (14)	high, full (1 exc)	late (1 exc)	late	TRUE	Zero, soon	0 (1 exc)	Zero <i>(1 exc)</i>	zero (1 exc)	soon
Stairway To Heaven	Quantum Steps (23) (17)	high, fair	Vp <b>0, early</b> (1 exc)	middle (2 exc)	FALSE	fair, long	0 – 3	Zero, few (2 exc)	zero, fair	fair (1 exc)
neuven	(17)	low, fair, high	middle	late	FALSE	fair, long	0 – 3	Zero, few, fair	zero, fair, high	soon (1 exc)
	Regularly Maintained (14) (11)	low, fair	V <sub>p</sub> 0, early (3 exc)	middle, late	FALSE	long vlong	>3	few, fair	high (1 exc)	soon, fair
	(3)	fair	middle	late	FALSE	fair, long	>3	few, fair	high (1 exc)	soon
Scared To Fall Asleep	Siesta (10)	fair (2 exc)	$V_p0$ , early	late	FALSE	vlong (1 exc)	<b>0 – 3</b> (2 exc)	Zero, few	zero, fair, high	soon
Again	Smoking Funnel (7)	fair (1 exc)	middle	middle	FALSE	fair	>3	fair, large	fair, high	fair

### Disjointness

- Essential Disjointness: are projects "occupying" solely an area of the space of actual values?
- To a very large extent yes! the patterns are focused in a specific area of the domain space and disjoint from the others.

#### • Exceptions:

- a couple of Siesta projects overlapping with Regularly Curated projects of similar definition,
- the Quantum Steps and Regularly Curated patterns:
  - although disjoint...
  - ... they are the only patterns practically separated by change rate in their growth period ...
  - ... and span a large area of the domain space of values
- A decision tree is also produced to verify essential Disjointness (with only 4 out of 151 projects that would have been erroneously classified)







- Genuine?
- Cohesive & disjoint?
- Generalizable?
- Complete?

### Visual inspection

Pattern cohesion refers to the internal homogeneity of the projects that pertain to each pattern.

Our first attempt towards producing cohesive patterns was via visual inspection.

The visual inspection of the schema evolution progress of the different patterns is quite revealing on the similarity between the members of each pattern.

# Cohesion: Just a few exceptions exist

- Genuine?
- Cohesive & disjoint?
- Generalizable?
- Complete?

- Two projects classified as Sigmoid violate the "middle-born" part of the definition by being born early.
- A Late Riser project reaches the top band in middle life, violating the requirement of late attainment of topband.
- Siesta has 2 projects exceeding the 0-3 months growth activity in the end, and a project that reaches growth just 'long' after schema birth (and not 'very long').
- A Quantum Step project reaches top late rather than middle.

Pattern	#prjs	Exceptions	Overlaps
Flatliner	23	-	_
Radical Sign	41	-	-
Sigmoid	19	2	-
Late Riser	14	1	-
Quantum Steps	23	2	-
<b>Regularly Curated</b>	14	-	-
Smoking Funnel	7	—	-
Siesta	10	3	-

See report in the accompanying web site of the paper

### Cohesion: quantitative metrics

- Genuine?
- Cohesive & disjoint?
- Generalizable?
- Complete?
- We quantized each project's time series to a vector of 20 measurements, one for each interval of 5% of time (i.e., at 0%, 5%, 10%, . . . of time), and computed the centroid for the corpus of each pattern.
- We assessed the SSE and the Mean Distance to Centroid per pattern

Pattern	#prjs	SSE	MDC
11_FlatLiner	23	0.18	0.06
12_RadicalSign	41	23.25	0.70
13_Sigmoid	19	31.76	1.26
14_LateRiser	14	9.33	0.78
21_QuantumSteps	23	36.20	1.16
22_RegularlyCurated	14	12.64	0.90
31_SmokingFunnel	7	5.00	0.83
32_Siesta	10	9.15	0.91

$$SSE(C_i) = \sum_{x \in C_i} dist(x, c_i)^2$$
$$MDC(C_i) = \frac{1}{m_i} \sum_{x \in C_i} dist(x, c_i)$$

1

- C<sub>i</sub> refers to cluster *i*, with a center
   c<sub>i</sub> and m<sub>i</sub> data points
- x is a data point assigned to C<sub>i</sub>
- dist(a,b) is a distance function
   between two data points, which,
   unless stated otherwise, we
   assume to be the Euclidean
   distance

# Threats to validity: External Validity

- Genuine?
- Cohesive & disjoint?
- Generalizable?
- Complete?
- External validity: we argue that the larger dataset is a very good representative of significant FOSS projects
  - GitHub is the main public repository for FoSS prj's.
  - We applied the filter of more than one contributor, more than 0 stars and non-forking
  - Subsequently, filtered out tests, examples & demos
  - Project domains include Content Management Systems, IoT Management on the cloud, Task Management Systems for O/S's, Messaging Platforms, Systems for the management of Scientific Data, Web on-line stores, On-line Charging Systems (OCS)...
- Limitations:
  - Multi-vendor DDL: covered only one vendor
  - Non-sql schemata, non .sql suffixes, multi-file DDL, incremental definitions of DDL



# Threats to validity: External Validity

- External validity:
  - The corpus of 151 significant projects is substantial and (b) suitable: with more than 12 months duration the projects hade ample time for the patterns to appear.
  - The patterns produced are **cohesive**, **disjoint** and **not overfitted** to the data set used.
  - This means we can generalize our results to the broader scope of schema evolution for FOSS

# Threats to validity: Experimental Validity

- We tested our extraction scripts with OpenCart, the largest of our studied projects for which we had a previous past extraction of its history, in 2016.
  - almost identical result, as only one commit out of 412 was missing from the GitHub history we extracted.
- 100% match for manual test of the histories of the retrieved files for a random sample of 50 cases.
- 100% match for removed projects from GitHub at the time of the cloning via a sample of 7 of them.
- Concerning our own software, we did extensive checks to our metrics computation tools.

# Threats to validity: Experimental Validity

- Experimental validity:
  - We have checked our change detection and chartgenerating software via tests and code reviews & trust it.
  - We have iteratively rechecked our manual allocation of projects to patterns

#### Is what se see ... complete?

- Genuine?
- Cohesive & disjoint?
- Generalizable?
- Complete?
- We cannot exclude the possibility that other patterns of change do exist
- Still, the overall situation hints that the existence of other patterns is rather unexpected:
- A large part of the space of possible value combinations is already covered.
- Second, although there do exist value combinations several of them are unattainable
  - e.g., a project with late schema birth, is obligatorily restricted to have a late top-band attainment and a short tail.

# Patterns and their characteristics

Explicit relationship of patterns to schema evolution quantitative characteristics



#### Project Update Period (PUP) in months

It is extremely revealing that almost all patterns have median values of the Project Update Period that are very close.

This means that the behavior of the schema evolution does NOT relate to the duration of the project – i.e., there is no "linear" (or other) relationship between how long the project lasted, wrt how its schema evolved!

The only half-exception seems to be Regularly Curated projects, but only wrt the median – the boxplot is very much included within the box plots of Late Risers and Smoking Funnels



#### Schema Update Period (# months)

Very expectedly, there is a progressive shift of the median value and the box plot:

- Almost frozen patterns obviously have very small numbers – o(3 months)
- Radical Sign and Quantum have median SUP periods around 11 and 12 months
- Smoking Funnel and Siesta around 31 months
- Regularly Curated: around 47.5 months

80



#### Schema Size at Birth (attributes)

- Almost all patterns have similar volumes of newborn schemata.
- Most patterns have their median value found around 22 attributes ...
- .. with the notable exception of Regularly Curated (62 attributes) and Smoking Funnel (79 attributes).
- Tables are very similar (3-5 tables at start for most patterns, vs 9.5 and 12 for the 2 exceptions)

1400

1200

#### **Total Activity**



Total Schema Activity (as well as Total Schema Expansion and Total Schema Maintenance which typically follow the pattern of Total Activity) is rather indifferent to the patterns.

Change values are small ...

.. except for Smoking Funnel and Regular Curation (the two smallest patterns in terms of population), who have larger values.

Progressive shift from

- a) top four, "almost frozen" patterns,
- b) Siesta and Quantum steps with somewhat higher medians, and,
- c) Smoking Funnel and Regular Curation which come with orders-ofmagnitude higher values.

However, this monotonicity does not change the fact that the overlap of the fairly narrow IQR's, at low values is significant.

### Total Activity

- Why do we see this behavior? It is important to highlight that the behavior towards schema evolution is not obligatorily in sync with the behavior towards source code evolution. For example, even in really frozen patterns, like, e.g., the populous Radical Sign, we can see vast ranges of PUP values. Similarly, for the schema size at birth.
- Why Smoking Funnel and Regularly Curated projects are different? These two include the most active projects in terms of total amount of change.

# Summary of how patterns differ wrt evolution characteristics

- Almost all patterns have median values of the Project Update Period that are very close, i.e., their evolution does NOT relate to the duration of the project
- Almost all patterns have similar volumes of new-born schemata, except for Regularly Curated and Smoking Funnel that have significantly larger mean values of birth volume
- Total Schema Activity (as well as Total Schema Expansion and Total Schema Maintenance which typically follow it) is rather indifferent to the patterns, with small change values except for Regularly Curated and Smoking Funnel who have larger values than the rest.
- Thus, Smoking Funnel and Regularly Curated projects start bigger and contain higher amounts of schema evolution activity than the rest of the projects who start small and typically show lower values of change.

### Predicting Patterns

Can we predict which pattern a project might follow, given some early characteristics?

# Predicting schema evolution is an extremely difficult problem

- To the best of our knowledge, there is no attempt to the problem so far for two reasons.
  - To a very large extent, the evolution depends heavily on the particularities of the development and curation team, along with the idiosyncrasies of the project itself. Thus, predictions would require a very detailed charting of both project and anthropocentric characteristics.
  - A really large corpus of schema histories is needed to come up with some broad approximation from a statistical point of view.
- We make such a preliminary attempt in this paper, to give a broad overview of how the point of schema birth and the future behavior of a project in terms of its pattern of schema evolution re-late.
- Several attempts towards further refinements (e.g., in the form of decision trees) have not provided better insights (OK, "not yet" ☺).

## Predictions are always hard, esp., for the future

Research Question: "Assume a curator, or an external assessor, who extracts (e.g., via git log and out tool set) the history of changes of a software project and its relational database. Can the curator make an educated guess on the future of how the schema will evolve?"

	151	100%	52	100%	38	100%	13	100%	48	100%
Siesta	10	7%	6	12%	3	8%	1	8%		
Smoking Funnel	7	5%							7	15%
Regularly Curated	14	9%	3	6%	4	11%	3	23%	4	8%
Quantum Steps	23	15%	4	8%	11	29%	2	15%	6	13%
Late Risers	14	9%							14	29%
Sigmoid	19	13%			1	3%	2	15%	16	33%
RadicalSign	41	27%	16	31%	19	50%	5	38%	1	2%
Flatliners	23	15%	23	44%						
_	#prj's	prob (%)	#prj's	prob (%)	#prj's	prob (%)	#prj's	prob (%)	#prj's	prob (%)
	OVER	ALL	Born	M0	Born	[M1 M6]	Born [M	I7M12]	Not Born	till M12

Assuming that the schema is born during M0 of the project life, the probability that the schema will be

- completely frozen is 75% (!), as a flatliner or a radical sign.
- under steady regular curation via the Quantum Steps or Regular Curation patterns, is 14%.
- 1/3 of the corpus projects have been born at month zero.

#### Predicting schema evolution

	OVERALL		Born M0		Born [M1 M6]		Born [M7M12]		Not Born till M12	
	#prj's	prob (%)	#prj's	prob (%)	#prj's	prob (%)	#prj's	prob (%)	#prj's	prob (%)
Flatliners	23	15%	23	44%						
RadicalSign	41	27%	16	31%	19	50%	5	38%	1	2%
Sigmoid	19	13%			1	3%	2	15%	16	33%
Late Risers	14	9%							14	29%
Quantum Steps	23	15%	4	8%	11	29%	2	15%	6	13%
Regularly Curated	14	9%	3	6%	4	11%	3	23%	4	8%
Smoking Funnel	7	5%							7	15%
Siesta	10	7%	6	12%	3	8%	1	8%		
	151	100%	52	100%	38	100%	13	100%	48	100%

- Assuming the schema is born in the next six months, there is ...
- ... a 53% chance that the schema will follow a sharp, focused evolution, and,
- ... a 40% chance it will follow a regularly curated life.
- 90 projects out of 151, i.e., 60% of the corpus, have already been born in the first six months of the projects' lives.

#### Predicting schema evolution

	OVERALL		Born M0		Born	[M1 M6]	Born [M7M12]		Not Born till M12	
_	#prj's	prob (%)	#prj's	prob (%)	#prj's	prob (%)	#prj's	prob (%)	#prj's	prob (%)
Flatliners	23	15%	23	44%						
RadicalSign	41	27%	16	31%	19	50%	5	38%	1	2%
Sigmoid	19	13%			1	3%	2	15%	16	33%
Late Risers	14	9%							14	29%
Quantum Steps	23	15%	4	8%	11	29%	2	15%	6	13%
Regularly Curated	14	9%	3	6%	4	11%	3	23%	4	8%
Smoking Funnel	7	5%							7	15%
Siesta	10	7%	6	12%	3	8%	1	8%		
	151	100%	52	100%	38	100%	13	100%	48	100%
- Assuming the schema is born in months 7 12, there is ...
- a drop in the number of projects born (13) ...
- ... but with similar distribution of probabilities for the families.
- The only significant change is that the probability of a more active regular curation rises to 23%.
- 2/3 of the corpus' projects have the schema being born within the first year.

## Predicting schema evolution

	OVER	ALL	Borr	n M0	Born	[M1 M6]	Born [N	17M12]	Not Born	till M12
	#prj's	prob (%)	#prj's	prob (%)	#prj's	prob (%)	#prj's	prob (%)	#prj's	prob (%)
Flatliners	23	15%	23	44%						
RadicalSign	41	27%	16	31%	19	50%	5	38%	1	2%
Sigmoid	19	13%			1	3%	2	15%	16	33%
Late Risers	14	9%							14	29%
Quantum Steps	23	15%	4	8%	11	29%	2	15%	6	13%
Regularly Curated	14	9%	3	6%	4	11%	3	23%	4	8%
Smoking Funnel	7	5%							7	15%
Siesta	10	7%	6	12%	3	8%	1	8%		
	151	100%	52	100%	38	100%	13	100%	48	100%

- As the point of schema birth exceeds the first year, there is
- ... a 64% chance of a sharp focused change,
- ... a 21% chance for a regular maintenance, and
- ... a 15% chance for a smoking funnel behavior.

## Predicting schema evolution

	OVER	ALL	Born	n M0	Born	[M1 M6]	Born [N	И7M12]	Not Born	till M12
	#prj's	prob (%)	#prj's	prob (%)	#prj's	prob (%)	#prj's	prob (%)	#prj's	prob (%)
Flatliners	23	15%	23	44%						
RadicalSign	41	27%	16	31%	19	50%	5	38%	1	2%
Sigmoid	19	13%			1	3%	2	15%	16	33%
Late Risers	14	9%							14	29%
Quantum Steps	23	15%	4	8%	11	29%	2	15%	6	13%
Regularly Curated	14	9%	3	6%	4	11%	3	23%	4	8%
Smoking Funnel	7	5%							7	15%
Siesta	10	7%	6	12%	3	8%	1	8%		
	151	100%	52	100%	38	100%	13	100%	48	100%

# Mixture of change types

What does schema evolution consist of?

See report in the accompanying web site of the paper



	Expansion	Maintenance	#prj	%prjs of pattern
11_FlatLiner	0%	0%	11	48%
23	0%	100%	5	22%
	(0%-40%)	(60%-100%)	0	0%
	middle	level	1	4%
	(60%-100%)	(0%-40%)	4	17%
	100%	0%	2	9%
12_RadicalSign	0%	0%	6	15%
41	0%	100%	1	2%
	(0%-40%)	(60%-100%)	6	15%
	middle	level	9	22%
	(60%-100%)	(0%-40%)	9	22%
	100%	0%	10	24%
13_Sigmoid	0%	0%	8	42%
19	0%	100%	1	5%
	(0%-40%)	(60%-100%)	1	5%
	middle	level	2	11%
	(60%-100%)	(0%-40%)	3	16%
	100%	0%	4	21%
14_LateRiser	0%	0%	2	14%
14	0%	100%	3	21%
	(0%-40%)	(60%-100%)	0	0%
	middle	level	2	14%
	(60%-100%)	(0%-40%)	1	/%
24. Question Otoma	100%	0%	6	43%
21_Quantumsteps	0%	0%	0	0%
23	(0% 40%)	100%	4	17%
	(0%-40%) middlo	(00%-100%)	4	1/%
	(60% 100%)	(0% 40%)	э 0	1370
	100%	(0%-40%)	0	170
22 RegularlyCurated	100%	0%	4	17%
22_RegulariyCurateu	0%	100%	0	0%
14	(0%-40%)	(60%-100%)	0	0%
	middle	(0070-10070)	7	50%
	(60%-100%)	(0%-40%)	7	50%
	100%	0%	0	0%
31 SmokingFunnel	0%	0%	0	0%
7	0%	100%	0	0%
	(0%-40%)	(60%-100%)	0	0%
	middle	level	5	71%
	(60%-100%)	(0%-40%)	2	29%
	100%	0%	0	0%
32 Siesta	0%	0%	0	0%
10	0%	100%	1	10%
	(0%-40%)	(60%-100%)	2	20%
	middle	. level /	4	40%
	(60%-100%)	(0%-40%)	3	30%
	100%	. 0%	0	0%

**Flatliners**. Mostly zero change. -- the rest, equally split, mostly monothematically, in a maintenance only, and, mostly expansion.

Expansion is mostly achieved via attribute injections. Maintenance is mostly manifested via type updates.

Radical Sign. Archetype "build; few changes; freeze" – most populous of all.
Radical sign change is predominantly achieved via table insertions and deletions.
1/3 of the projects are maintenance oriented; the rest, towards expansion, too often monothematic.

**Sigmoid**. 8 out of 19 projects have zero change; the rest have mostly low change, mostly monothematic (5 /12). Projects are slightly inclined towards expansion, mostly via attribute injections. There is no maintenance really, except for two large projects with too many table deletions.

**Late Riser**. projects with small change (with 1 exception), heavily monothematic (9/12). The pattern is slightly biased towards expansion. Maintenance: 4 projects inclined to maintenance, mostly via type updates. The rest: (small volume) expansion-only projects, mostly via attribute injections.

**Quantum Steps**. 4/23 projects that are maintenance-only, with change manifested via type updates. When maintenance dominates: mostly via update-type. Mixed nature projects: maintenance is mostly achieved via deletion of entire tables. When expansion is predominant (12/23): via new table births.

**Regularly Curated**. Change rises to significant volumes. Heavily mixed (only 2 rows in the table): almost all categories (exc. PK) appearing in almost all projects. The mix of change is mostly balanced, slightly in favor of expansions. Expansion is mostly achieved via new tables. Injections clearly exist but significantly less than attributes been born with new tables. Maintenance is mostly achieved via table deletions, with (i) attributes being ejected and (ii) data types being updated constituting a second tier of types in terms of popularity.

Smoking Funnel. Similar to the one of Regularly Curated projects.

Siesta. Heavily mixed, with each project having a large number of change types
 (at least 3). In contrast to the other categories, siesta projects are slightly in favor
 of maintenance. Maintenance is split half and half between (a) table deletions
 and (b) data type updates. Expansion is mostly achieved via table insertions.

## Summary of Schema Change Types

- With the exception of few radical sign projects that are oriented towards maintenance, the "Be Quick or Be Dead" family involves small change, frequently being zero, and an inclination towards expansion. Due to the small volume of change, the patterns of the family are frequently monothematic in their internal breakdown of change types.
- In contrast, the rest of the time-related patterns come with higher volumes of change, which is also related to a variety of change types. Both expansion and maintenance are performed with the granule of change being mostly the entire table (being inserted or deleted), rather than the internal restructuring of existing tables.

Lessons Learned and Open Issues

... and why it matters...

## Main contribution

- For the first time in the literature we know how schema evolution unfolds over time
- We have extracted ...
- ... 8 Patterns of Change in Time...
- ... organized in 3 Families
- We verified the results of pattern extraction wrt common sense, generalization, disjointness, cohesion, completeness

	<ul> <li>Count of priss</li> </ul>
□ 1_BeQuickOrBeDead	97
11_FlatLiner	23
12_RadicalSign	41
13_Sigmoid	19
14_LateRiser	14
2_StairwayToHeaven	37
21_QuantumSteps	23
22_RegularlyCurated	14
□ 3_ScaredToFallAsleepAgain	n 17
31_SmokingFunnel	7
32_Siesta	10
Grand Total	151



## Aversion to change

- Almost 2/3 of the corpus, 97/151 evolve with very focused change very close to the point of schema birth the only difference of the involved patterns is when schema birth takes place.
- This is in sync with our recent previous findings over a very large corpus of projects, where 70% of 327 projects investigated showed very little – if any- signs of change.



 This is in contrast with all the past research till the 20's, involving small studies, of few, carefully picked prj's that actually showed any change whatsoever

## Still, change exists!

 Almost 25% (37/151) projects evolve with regular change steps, distributed across time, progressively climbing to the top-band over a long period of time.





 There is even late change! 17/151 projects come with change not focused in a single point, happening towards the end of the lifetime of the project.





## Summary of simple stats

- Birth is mostly done early
  - Two thirds of the projects (105 projects) see schema birth at V0 or before 25% of the PUP.
  - 74 schemata (half the corpus) are born in the first 10% of time.
- Change is mostly sharp in time:
  - 115 projects, 76% of the population had no more than 1 month of activity from schema birth to top.
  - 88 / 151 projects (58%) had a single vault, i.e., an interval from schema birth to top-band that was less than 10% of the PUP, with 62 of them in zero time.

## Summary of Schema Change Types

- With the exception of few radical sign projects that are oriented towards maintenance, the "Be Quick or Be Dead" family involves small change, frequently being zero, and an inclination towards expansion. Due to the small volume of change, the patterns of the family are frequently monothematic in their internal breakdown of change types.
- In contrast, the rest of the time-related patterns come with higher volumes of change, which is also related to a variety of change types. Both expansion and maintenance are performed with the granule of change being mostly the entire table (being inserted or deleted), rather than the internal restructuring of existing tables.

# Summary of how patterns differ wrt evolution characteristics

- Almost all patterns have median values of the Project Update Period that are very close, i.e., their evolution does NOT relate to the duration of the project
- Almost all patterns have similar volumes of new-born schemata, except for Regularly Curated and Smoking Funnel that have significantly larger mean values of birth volume
- Total Schema Activity (as well as Total Schema Expansion and Total Schema Maintenance which typically follow it ) is rather indifferent to the patterns, with small change values except for Regularly Curated and Smoking Funnel who have larger values than the rest.
- Thus, Smoking Funnel and Regularly Curated projects start bigger and contain higher amounts of schema evolution activity than the rest of the projects who start small and typically show lower values of change.

When are schemata born within a project's life?

- 34% of the schemata are born in M0
- 60% of the schemata are born in the first 6 months
- 68% of the schemata are born in the first 12 months

	OVER	ALL	Born	M0	Born	[M1 M6]	Born [N	M7M12]	Not Bor	n till M12
_	#prj's	prob (%)	#prj's	prob (%)	#prj's	prob (%)	#prj's	prob (%)	#prj's	prob (%)
Flatliners	23	15%	23	44%						
RadicalSign	41	27%	16	31%	19	50%	5	38%	1	2%
Sigmoid	19	13%			1	3%	2	15%	16	33%
Late Risers	14	9%							14	29%
Quantum Steps	23	15%	4	8%	11	29%	2	15%	6	13%
Regularly Curated	14	9%	3	6%	4	11%	3	23%	4	8%
Smoking Funnel	7	5%							7	15%
Siesta	10	7%	6	12%	3	8%	1	8%		
	151	100%	52	100%	38	100%	13	100%	48	100%

# What's the chance of the schema freezing depending on when it was born?

Born in	% chance to end up as Be quick or be dead	% chance to end up with some kind of regular change
M0	75%	25%
[M01-M06]	53%	47%
[M07 – M12]	53%	47%
>M12	64%	36%

	OVER	ALL	Born	M0	Born	[M1 M6]	Born [N	//7M12]	Not Bori	n till M12
_	#prj's	prob (%)	#prj's	prob (%)	#prj's	prob (%)	#prj's	prob (%)	#prj's	prob (%)
Flatliners	23	15%	23	44%						
RadicalSign	41	27%	16	31%	19	50%	5	38%	1	2%
Sigmoid	19	13%			1	3%	2	15%	16	33%
Late Risers	14	9%							14	29%
Quantum Steps	23	15%	4	8%	11	29%	2	15%	6	13%
Regularly Curated	14	9%	3	6%	4	11%	3	23%	4	8%
Smoking Funnel	7	5%							7	15%
Siesta	10	7%	6	12%	3	8%	1	8%		
	151	100%	52	100%	38	100%	13	100%	48	100%

## Lessons learned and who cares

- The "freeze-and-build" aversion to change is indeed present at large and concerns 2/3 of the corpus (Be Quick or Be Dead family).
  - For researchers: we need to reflect on our schema design and software development premises.
  - For Developers: To battle aversion to change, development teams should chart the mapping of schema to source code in order to be able to adapt them in sync.
- The rest of the projects evolve differently in time: rare but regular change; densely regular change; and, surprisingly, late change too. Thus, we cannot rely solely on the premise of aversion to change.
- We also have insights on how likely it is to encounter change or rigidity in the life of a schema.
  - Project managers can reserve time up-front for handling change and its impact.
- Sadly atypically for our data engineering community, we contribute with the enrichment of our knowledge —in a principled manner— on how the artifacts that we invent (here: relational databases) are actually used in practice. Part of our contribution is the topic itself, along with the nomenclature, measures and methodological principles used.

## With an eye to the future

- Methodologically, the paper opens a road for future research on other kinds of schemata (e.g., do the same for Nosql schemata)
- Solid foundations for the **prediction of future behavior** on the basis of a meaningful model.
- Obtain and study schema histories from proprietary schemata (for the last 50 years of the database discipline, this has been practically impossible).
- (Even more importantly) **Developing educational material and practical exercises for our students**, in order to train them on the practical aspects of the topic is another important road for the future.

## Thank you!

## https://github.com/DAINTINESS-Group/

### **Everything is online!**

My group's git page

https://github.com/DAINTINESS-Group/

has links to **Data sets** 

https://github.com/DAINTINESS-Group/Schema\_Evolution\_Datasets/tree/mast er/SchemaEvolutionDatasets2020

and Code

- ... for computing differences (Hecate)
- ... visualizing schema lives (Plutarch Par. Lives)
- ... visualizing the structure of FK's (Parmenidian Truth)
- ... handling the impact of evolution (Hecataeus)

daintiness bas Intervie Internation Comparents Der uf Comp Science & Engineering Entervietige Elements	DAta INTensive Information DAta INTensive Information EcoSystemS Group ⊘ Ioannina, Greece	n EcoSystemS Group a, Univ. Ioannina, Hellas
📮 Repositorie	s 12 🔗 Packages 🔗 People 4 🛄	] Projects
Q Find a repos	itory	Type 👻 Language 🕶

This is a project for the monitoring of the evolution of the parallel histories of entities that evolve in parallel. This is a new, that builds upon the previous Plutarch Parallel Lives (attn to the underscores at the name) that visualized the schema evolution of the tables of a relational schema.

visualization iava

Java 🖞 0 🏠 0 🕕 0 🚺 0 Updated 7 days ago

Schema_Evolution_Datasets	
Forked from giskou/EvolutionDatasets Collections of schema histories, to be studied for their schema evolution.	^^
sql relational-databases schema-evolution	
■ PLSQL ¥ 2 ☆ 3 ① 0 \$ 0 Updated on Jan 23	

### Hecate Forked from giskou/Hecate Diff visualization between 2 SQL schemas relational-databases schema-evolution

▲ AGPL-3.0 💡 6 ☆ 1 ① 0 🏌 0 Updated on Dec 10, 2020

#### ParmenidianTruth

Visualizes the story of a database's schema as a pptx presentation ● Java 学0 ☆0 ① 0 🎵 0 Updated on Oct 10, 2018

To probe further (code, data, details, presentations, ...) www.cs.uoi.gr/~pvassil/projects/schemaBiographies/