

AN ABSTRACT OF THE DISSERTATION OF

Parisa Ataei for the degree of Doctor of Philosophy in Computer Science presented on
June 20, 2021.

Title: Theory and Implementation of a Variational Database Management System

Abstract approved: _____

Eric Walkingshaw

In this thesis I present the variational database management system, a formal framework and its implementation for representing variation in relational databases and managing variational information needs. A variational database is intended to support any kind of variation in a database. Specific kinds of variation in databases have already been studied and are well-supported, for example, schema evolution systems address the variation of a database's schema over time and data integration systems address variation caused by accessing data from multiple data sources simultaneously. However, many other kinds of variation in databases arise in practice, and different kinds of variation often interact, but these scenarios are not well-supported by the existing work. For example, neither the schema evolution systems nor the database integration systems can address variation that arises when data sources combined in one database evolve over time.

This thesis collects a large amount of work: It defines the variational database framework and the syntax and [specific kind of] semantics of the variational relational algebra, a query language for variational databases. It also defines the requirements of a generic variational database framework that makes the framework expressive enough to encode any kind of variation in databases. Additionally, it [shows/proves] that the introduced framework satisfies all these needs. It presents two use cases of the variational database framework that are based on existing data sets and scenarios that are partially supported by existing techniques. It presents the variational database management system which is the implementation of variational databases and variational relational algebra as an abstract layer written in Haskell on top of a traditional RDBMS. It also presents several

theoretical results related to the framework and query language, such as syntax-based equivalence rules that preserve the semantics of a query, a type system for ensuring that a variational query is well formed with respect to the underlying variational schema, and a confluence property of the variational relational algebra type system and semantics with respect to the relational algebra type system and semantics.

©Copyright by Parisa Ataei
June 20, 2021
All Rights Reserved

Theory and Implementation of a Variational Database Management System

by

Parisa Ataei

A DISSERTATION

submitted to

Oregon State University

in partial fulfillment of
the requirements for the
degree of

Doctor of Philosophy

Presented June 20, 2021
Commencement June 2021

Doctor of Philosophy dissertation of Parisa Ataei presented on June 20, 2021.

APPROVED:

Major Professor, representing Computer Science

Director of the School of Electrical Engineering and Computer Science

Dean of the Graduate School

I understand that my dissertation will become part of the permanent collection of Oregon State University libraries. My signature below authorizes release of my dissertation to any reader upon request.

Parisa Ataei, Author

ACKNOWLEDGEMENTS

[Eric. Committee. jeff. abu. parents. friends.]

LIST OF ALGORITHMS

Algorithm

Page

Chapter 1: Introduction

Managing variation in databases is a perennial problem in database literature and appears in different forms and contexts [?, ?, ?, ?, ?] and it is inevitable [?]. Variation in databases mainly arises when multiple database instances conceptually represent the same database, but, differ in their schema, content, or constraints. Specific kinds of variation in databases have been addressed by context-specific solutions, such as schema evolution [?, ?, ?, ?, ?], data integration [?], and database versioning [?, ?]. However, there are no generic solution that addresses any kind of variation in databases. We motivate the need for a generic solution to variation in databases in Section ??.

The major contribution of this thesis is the *variational database* framework, a generic relational database framework that explicitly accounts for database changes (variation), and *variational relational algebra*, a data manipulation language for our framework that allows for information extraction from a variational database. The variational framework is generic because it can encode any kind of variation in databases. Additionally and more importantly, it is designed such that it satisfies any information need that a user may have in a variational database scenario. Based on information needs in a variational database scenario, we define requirements of a variational database framework in Section ?? and throughout the thesis, we show that our framework satisfies these requirements.

In addition to a formal description of the variational database framework and variational relational algebra and some theoretical results, this thesis distributes and presents two variational data sets (including both a variational database and a set of queries) as well as a *variational database management system* that implements the variational database framework as an abstraction layer on top of a traditional relational database management system in Haskell. Section ?? enumerates the specific contributions in the context of an outline of the structure of the rest of the thesis.

1.1 Motivation and Impact

Managing variation in databases is a perennial problem in database literature and appears in different forms and contexts [?, ?, ?, ?, ?]. Variation in databases mainly arises when multiple database instances conceptually represent the same database, but, differ in their schema or content. The variation in schema and/or content occurs in two *dimensions*: time and space. Variation in space refers to different variants of database that coexist in parallel while variation in time refers to the evolution of database, similar to variation observed in software [?]. Note that variation in a database can occur due to both dimensions at the same time.

Existing work on variation in databases addresses specific kinds of variation in a context and proposes solutions specific to the context of the variation, such as schema evolution [?, ?, ?, ?, ?], data integration [?], and database versioning [?, ?]. Unfortunately, some of these tools do not address all their user’s needs. Furthermore, they are all *variation-unaware*, i.e., they dismiss that they are addressing a specific kind of a bigger problem. Thus, not only they cannot address a new kind of database variation but they also cannot address the interaction of different kinds of database variation since they assume that each kind of variation is *isolated* from another kinds. This costs database researchers to develop a new system for every individual kind of data variation and forces developers and DBAs to use manually unsafe workaround.

For example, schema evolution is a kind of schematic variation in databases over time that is well-supported [?, ?, ?, ?, ?]. Changes applied to the schema over time are *variation* in the database and every time the database evolves, a new *variant* is generated. Current solutions addressing schema evolution dismiss that it is a kind of variation, thus, they *simulate* the effect of variation by relying on temporal nature of schema evolution and using timestamps [?, ?, ?, ?] or keeping an external file of time-line history of changes applied to the database [?].

Unlike schema evolution, some kinds of database variation are partially supported. For example, while developing software for multiple clients simultaneously, an approach called software product line (SPL) [?], a different database schema is required for each client due to client’s different business requirements and environments [?]. SPL researchers have developed encodings of data models that allow for arbitrary variation by annotating different elements of the model with features from the SPL [?, ?, ?]. Thus,

they can generate a database schema *variant* for each software *variant* requested by a client and generated by the SPL. However, these solutions address *only* variation in the data model but do not extend it to the level of the data or queries. The lack of variation support in queries leads to unsafe techniques such as encoding different variants of query through string munging, while the lack of variation support in data precludes testing with multiple variants of a database at once.

The situation exacerbates even more when two kinds of variation interact and create a new kind of database variation: the evolution of the database used in development of software by an SPL. This is due to the evolution of software and its artifacts; an inevitable phenomena [?]. This is where even previous context-specific solutions like schema evolution tools fail because the context-specific approaches assume that the kind of variation that they address is completely isolated from other kinds of variation. We motivate this case through an example in Section ???. Note that new kinds of variation may arise that are not necessarily the interaction of already-addressed variation in database. We discuss such a scenario in Section ???.

As we have shown, variation in databases is abundant and inexorable [?]; impacts DBAs, data analysts, and developers significantly [?]; and appears in different contexts. Current methods are all extremely tailored to a specific context. Consequently, they fail to address the interaction of their specific kind of database variation with other kinds. Hence, the challenge becomes defining a variation-aware database that can model different kinds of variation in different contexts such that it satisfies different specialists' needs at different stages, e.g., development, information extraction, deployment, and testing.

1.1.1 Motivating Example

In this section, we motivate the interaction of two kinds of variation in databases resulting in a new kind: database-backed software produced by an SPL and schema evolution. An SPL uses a set of boolean variables called *features* to indicate the functionalities that each software variant requires. Consider an SPL that generates management software for companies. It has a feature *edu* that indicates whether a company provides educational means such as courses for its employees. Software variants in which *edu* is disabled (i.e., *edu* = *false*) only provide basic functionalities while ones in which *edu* is enabled

Table 1.1: Schema variants of the employee database developed for multiple software variants by an SPL. Note that an educational database variant must contain a basic database variant too.

(a) Database schema variants for basic software variants.

Temporal Features	basic Database Schema Variants
V_1	<i>engineerpersonnel</i> (<i>empno</i> , <i>name</i> , <i>hiredate</i> , <i>title</i> , <i>deptname</i>) <i>otherpersonnel</i> (<i>empno</i> , <i>name</i> , <i>hiredate</i> , <i>title</i> , <i>deptname</i>) <i>job</i> (<i>title</i> , <i>salary</i>)
V_2	<i>empacct</i> (<i>empno</i> , <i>name</i> , <i>hiredate</i> , <i>title</i> , <i>deptname</i>) <i>job</i> (<i>title</i> , <i>salary</i>)
V_3	<i>empacct</i> (<i>empno</i> , <i>name</i> , <i>hiredate</i> , <i>title</i> , <i>deptno</i>) <i>job</i> (<i>title</i> , <i>salary</i>) <i>dept</i> (<i>deptname</i> , <i>deptno</i> , <i>managerno</i>) <i>empbio</i> (<i>empno</i> , <i>sex</i> , <i>birthdate</i>)
V_4	<i>empacct</i> (<i>empno</i> , <i>hiredate</i> , <i>title</i> , <i>deptno</i> , <i>std</i> , <i>instr</i>) <i>job</i> (<i>title</i> , <i>salary</i>) <i>dept</i> (<i>deptname</i> , <i>deptno</i> , <i>managerno</i>) <i>empbio</i> (<i>empno</i> , <i>sex</i> , <i>birthdate</i> , <i>name</i>)
V_5	<i>empacct</i> (<i>empno</i> , <i>hiredate</i> , <i>title</i> , <i>deptno</i> , <i>std</i> , <i>instr</i> , <i>salary</i>) <i>dept</i> (<i>deptname</i> , <i>deptno</i> , <i>managerno</i> , <i>stdnum</i> , <i>instrnum</i>) <i>empbio</i> (<i>empno</i> , <i>sex</i> , <i>birthdate</i> , <i>firstname</i> , <i>lastname</i>)

(b) Database schema variants for educational software variants.

Temporal Feature	educational Database Schema Variants
T_1	<i>course</i> (<i>coursename</i> , <i>teacherno</i>) <i>student</i> (<i>studentno</i> , <i>coursename</i>)
T_2	<i>course</i> (<i>courseno</i> , <i>coursename</i> , <i>teacherno</i>) <i>student</i> (<i>studentno</i> , <i>courseno</i>)
T_3	<i>course</i> (<i>courseno</i> , <i>coursename</i>) <i>teach</i> (<i>teacherno</i> , <i>courseno</i>) <i>student</i> (<i>studentno</i> , <i>courseno</i> , <i>grade</i>)
T_4	<i>ecourse</i> (<i>courseno</i> , <i>coursename</i>) <i>course</i> (<i>courseno</i> , <i>coursename</i> , <i>time</i> , <i>class</i>) <i>teach</i> (<i>teacherno</i> , <i>courseno</i>) <i>student</i> (<i>studentno</i> , <i>courseno</i> , <i>grade</i>)
T_5	<i>ecourse</i> (<i>courseno</i> , <i>coursename</i> , <i>deptno</i>) <i>course</i> (<i>courseno</i> , <i>coursename</i> , <i>time</i> , <i>class</i> , <i>deptno</i>) <i>teach</i> (<i>teacherno</i> , <i>courseno</i>) <i>take</i> (<i>studentno</i> , <i>courseno</i> , <i>grade</i>)

provide educational functionalities in addition to the basic ones. Thus, this SPL yields two types of variants: **basic** and **educational**.

Each variant of this SPL needs a database to store information about employees, but SPL features impact the database: While **basic** variants do not need to store any education-related records **educational** variants do. We visualize the impact of features on the schema variants in Table ?? : It has two schema types: **basic**, shown in Table ?? , and **educational**, shown in Table ?? . A **basic** schema variant contains only the schema in one of the cells in Table ?? while an **educational** schema variant consists of two sub-schemas: one from the **basic** and another from Table ?? . For example, the yellow highlighted cells include relation schemas for an **educational** schema variant.

Rows of Table ?? indicate the evolution of schema variants in time. To capture the evolution of the software and its database we add two disjoint sets of features which again are boolean variables. The temporal feature sets are disjoint to allow individual paces for the evolution of each type of schema. For example, yellow cells of Table ?? show a valid schema variant even though the **basic** and **educational** sub-schemas are not in the same row.

Now, consider the following scenario: In the initial design of the **basic** database, SPL DBAs settle on three tables *engineerpersonnel*, *otherpersonnel*, and *job*; shown in Table ?? and associated with feature V_1 . After some time, they decide to refactor the schema to remove redundant tables, thus, they combine the two relations *engineerpersonnel* and *otherpersonnel* into one, *empacct*; associated with feature V_2 . Since some clients' software relies on a previous design the two schemas have to coexist in parallel. Therefore, the existence (presence) of *engineerpersonnel* and *otherpersonnel* relations is *variational*, i.e., they only exist in the **basic** schema when $V_1 = \text{true}$. This scenario describes *component evolution*: developers update, refactor, and improve components including the database [?].

Now, consider the case where a client that previously requested a **basic** variant of the management software has recently added courses to educate its employees in specific subjects. Hence, an SPL developer needs to enable the *edu* feature for this client, forcing the adjustment of the schema variant to **educational**. This case describes *product evolution*: database evolution in SPL resulted from clients adding/removing features/-components [?].

The situation is further complicated since the **basic** and **educational** schemas are in-

terdependent: Consider the **basic** schema variant for feature V_4 . Attributes *std* and *instr* only exists in the *empacct* relation when *edu* = **true**, represented by a dash-underline, otherwise the *empacct* relation has only four attributes: *empno*, *hiredate*, *title*, and *deptno*. Hence, the presence of attributes *std* and *instr* in *empacct* relation is *variational*, i.e., they only exist in *empacct* relation when *edu* = **true**.

Our example demonstrates how different kinds of variation interact with each other, an indispensable consequence of modern software development. The described interaction is similar to a recent scenario we discussed with an industry contact in Section ??.

1.1.2 New Instance of Data Variation in Industry

New variational scenarios could appear, either from combination of other scenarios or even a new scenario could reveal itself. For example, the following is a scenario we recently discussed with an industry contact: A software company develops software for different networking companies and analyzes data from its clients to advise them accordingly. The company records information from each of its clients' networks in databases customized to the particular hardware, operating systems, etc. that each client uses. The company analysts need to query information from all clients who agreed to share their information, but the same information need will be represented differently for each client. This problem is essentially a combination of the SPL variation problem (the company develops and maintains many databases that vary in structure and content) and the data integration problem (querying over many databases that vary in structure and content). However, neither the existing solutions from the SPL community nor database integration address both sides of the problem. Currently the company manually maintains variant schemas and queries, but this does not take advantage of sharing and is a major maintenance challenge. With a database encoding that supports explicit variation in schemas, content, and queries, the company could maintain a single variational database that can be configured for each client, import shared data into a variational database, and write variational queries over the variational database to analyze the data, significantly reducing redundancy across clients.

1.1.3 Requirements of a Variation-Aware Database Framework

For a variation-aware framework to be expressive enough to encode any kind of variation in databases, it must satisfy some requirements. Thus, we define the requirements that make a database framework variational through studying different kinds of variation in databases. A variational artifact, including databases, encompasses multiple *variants* of the artifact and provides a way to distinguish between different variants that are all encoded in one place, the variational artifact. It also provides a way to get the variants from the variational artifact, we call this *configuration function*. These requirements that help distinguish a database framework that simulates the effect of variation compared to one that is variational are listed below:

- (R0) *All database variants must be accessible at a given time.* For example, in our motivating example, a company that started with V_1 of the **basic** schema evolves over time but its different branches adopt the new schema at different paces, thus, it requires access to all variants of the **basic** schema.
- (R1) *The query language must allow for querying multiple database variants simultaneously. Additionally, it must allow for filtering tuples to specific variants.* That is, the framework must provide a query language that allows users to query multiple database variants at the same time in addition to giving them the freedom to choose the variants that they want to query. For example, an SPL tester that is testing a piece of code for the not highlighted variants of the software in Table ?? needs to write queries that exclude the variants associated with yellow cells of Table ??.
- (R2) *Every piece of data must clearly state the variant it belongs to and this information must be kept throughout the entire framework.* Continuing the example of the SPL tester, they need to know the variant that some results belong to in order to be able to debug the software correctly and accordingly.
- (R3) *The variational database must provide a way for generating database and query variants.* For example, the SPL developers need to deploy the management software for each client, thus, they need to configure the database schema and its queries in the code for each software variant.

Throughout the thesis, we show how the proposed variational database framework satisfies these requirements via examples, proofs, and tests.

1.2 Contributions and Outline of this Thesis

[to be filled out when I have the chapters.]

The high level goal of this thesis is to emphasize the need for a variation-aware database framework and to present one such framework. Therefore, in addition to the formal definition of the framework and query language, it also provides variational data sets (including both the variational database and a set of queries) to illustrate the feasibility of the proposed framework. Furthermore, it illustrates various approaches to implement such a framework and compares their performance.

The rest of this section describes the structure of this thesis, enumerating the specific contribution that each chapter makes.

Chapter ?? (*Background*) introduces several concepts and terms that are the basis of this thesis. It describes types and how to interpret them. It explains relational databases with assumptions that are held throughout the thesis and relational algebra. It also describes various ways of incorporating variation into elements of a database.

[Chapter ??]

[Chapter ??]

[Chapter ??]

[Chapter ??]

Chapter ?? (*Related Work*) collects research related to different kinds of variation in databases and other related variational research.

Finally, Chapter ?? (*Conclusion*) briefly presents ...

Chapter 2: Background

The core of this thesis is injecting a new aspect to relational databases: *variation*. Thus, the goals of this chapter are twofold: first, to introduce how variation is encoded and represented in our variational database framework; second, to provide the reader with the concepts and notations used to build up the main contributions of this thesis, mainly relational databases and relational algebra.

Section ?? describes the database model of relational databases and the specification of the structured used to store the data. Section ?? describes SPJC relational algebra to query relational databases. Then, Section ?? defines our choice of encoding for the variation space used in the variational database as propositional formulas of boolean variables. Finally, we introduce main approaches used to incorporate variation into our variational database framework. Section ?? introduces variation into sets which forms the basis of the variational database framework and Section ?? describes the formula choice calculus used to incorporate variation into relational algebra.

2.1 Relational Databases

[add examples of tables.] [add figure of all definitions]

A relational database \underline{D} stores information in a structured manner by forcing data to conform to a *schema* \underline{S} that is a finite set $\{\underline{s}_1, \dots, \underline{s}_n\}$ of *relation schemas*. A relation schema is defined as $\underline{s} = r(\underline{a}_1, \dots, \underline{a}_k)$ where each $\underline{a}_i \in \underline{\mathbf{A}}$ is an *attribute* contained in a relation named r . Note that there is a total order $\leq_{\underline{\mathbf{A}}}$ on $\underline{\mathbf{A}}$ and every listed set of attributes is written according to $\leq_{\underline{\mathbf{A}}}$. For theoretical development, it suffices to use the same domain of values for all of the attributes. Thus, we fix a countably infinite set *domain dom*. A *value* and a *constant* are elements of *dom*.

The content of database \underline{D} is stored in the form of *tuples*. A tuple \underline{u} is a mapping between an ordered set of attributes and their values, i.e., $\underline{u} = (\underline{v}_1, \dots, \underline{v}_k)$ for the relation schema $r(\underline{a}_1, \dots, \underline{a}_k)$. Hence a *relation content*, \underline{U} , is a set of tuples $\{\underline{u}_1, \dots, \underline{u}_m\}$. $att(i)$ returns the attribute corresponding to index i in a tuple. A *table* $\underline{t} = (\underline{s}, \underline{U})$ is a pair of

relation schema and relation content. A *database instance*, $\underline{\mathcal{I}}$, of the database \underline{D} with the schema \underline{S} , is a set of tables $\{\underline{t}_1, \dots, \underline{t}_n\}$. For brevity, when it is clear from the context we refer to a database instance by *database*.

A relational database \underline{D} stores information in a structured manner by forcing data to conform to a *schema* \underline{S} that is a finite set $\{\underline{s}_1, \dots, \underline{s}_n\}$ of *relation schemas*. A relation schema is defined as $\underline{s} = r(\underline{a}_1, \dots, \underline{a}_k)$ where each \underline{a}_i is an *attribute* contained in a relation named r . $rel(\underline{a})$ returns the relation that contains the attribute. $type(\underline{a})$ returns the *type* of values associated with attribute \underline{a} .

The content of database \underline{D} is stored in the form of *tuples*. A tuple \underline{u} is a mapping between a list of relation schema attributes and their values, i.e., $\underline{u} = (\underline{v}_1, \dots, \underline{v}_k)$ for the relation schema $r(\underline{a}_1, \dots, \underline{a}_k)$. Hence a *relation content*, \underline{U} , is a set of tuples $\{\underline{u}_1, \dots, \underline{u}_m\}$. $att(\underline{v})$ returns the attribute the value corresponds to. A *table* \underline{t} is a pair of relation content and relation schema. A *database instance*, $\underline{\mathcal{I}}$, of the database \underline{D} with the schema \underline{S} , is a set of tables $\{\underline{t}_1, \dots, \underline{t}_n\}$. For brevity, when it is clear from the context we refer to a database instance by *database*.

2.2 The SPJR Relational Algebra

[relational algebra] [add syntax definition] [add type system] [add examples with tables]
[maybe add semantics later on]

We do not extend the notation of using underline for relational algebra operations. Instead, relational algebra operations are overloaded and they are used as both plain relational and variational operations. It is clear from the context when an operation is variational or not. We also extend relational algebra s.t. projection of an empty list from a query is valid and it returns an empty set. In fact, we denote such query by the *empty* query ε , thus, $\varepsilon = \pi_{\{\}} q$.

[the following is from prelim. revise.]

Figure ?? defines the syntax of relational algebra which allows users to query a relational database [?]. The first five constructs are adapted from relational algebra: A query may simply *reference* a relation \underline{r} in the schema. *Renaming* allows giving a name to an intermediate query to be referenced later. Note that \underline{r} is an overloaded symbol that indicates both a relation and a relation name. A *projection* enables selecting a subset of attributes from the results of a subquery, for example, $\pi_{\underline{a}_1} \underline{r}$ would return only

$$\underline{\theta} \in \underline{\Theta} \quad := \quad \text{true} \mid \text{false} \mid a \bullet k \mid a \bullet a \mid \neg \underline{\theta} \mid \underline{\theta} \vee \underline{\theta}$$

$$\begin{array}{ll} \underline{q} \in \underline{\mathbf{Q}} & := \underline{r} \quad \text{Relation reference} \\ & \mid \rho_{\underline{r}} \underline{q} \quad \text{Renaming} \\ & \mid \pi_{\underline{A}} \underline{q} \quad \text{Projection} \\ & \mid \sigma_{\underline{\theta}} \underline{q} \quad \text{Selection} \\ & \mid \underline{q} \times \underline{q} \quad \text{Cartesian product} \\ & \mid \underline{q} \bowtie_{\underline{\theta}} \underline{q} \quad \text{Join} \\ & \mid \underline{q} \circ \underline{q} \quad \text{Set operation} \end{array}$$

Figure 2.1: Syntax of relational algebra, where \bullet ranges over comparison operators ($<, \leq, =, \neq, >, \geq$), \circ over set operations (\cap, \cup), k over constant values, a over attribute names, and \underline{A} over lists of attributes. The syntactic category $\underline{\theta}$ is relational conditions, and \underline{q} is relational algebra terms.

attribute \underline{a}_1 from \underline{r} . A *selection* enables filtering the tuples returned by a subquery based on a given condition $\underline{\theta}$, for example, $\sigma_{\underline{a}_1 > 3} \underline{r}$ would return all tuples from \underline{r} where the value for \underline{a}_1 is greater than 3. A *Cartesian products* simply cross products every tuple from its left subquery with every tuple from its right subquery. The *join* operation joins two subqueries based on a condition and omitting its condition implies it is a natural join (i.e., join on the shared attribute of the two subqueries). For example, $\underline{r}_1 \bowtie_{\underline{a}_1 = \underline{a}_2} \underline{r}_2$ joins tuples from \underline{r}_1 and \underline{r}_2 where the attribute \underline{a}_1 from relation \underline{r}_1 is equal to attribute \underline{a}_2 from relation \underline{r}_2 . However, if we have $\underline{r}_1(\underline{a}_1, \underline{a}_3)$ and $\underline{r}_2(\underline{a}_1, \underline{a}_2)$ then $\underline{r}_1 \bowtie \underline{r}_2$ joins tuples from \underline{r}_1 and \underline{r}_2 where attribute \underline{a}_1 has the same value in \underline{r}_1 and \underline{r}_2 . Also, note that join is simply a syntactic sugar for selection of cross product, that is $\underline{q}_1 \bowtie_{\underline{\theta}} \underline{q}_2 = \sigma_{\underline{\theta}}(\underline{q}_1 \times \underline{q}_2)$. The set operations, union and intersection, require two subqueries to have the same set of attributes and simply apply the operation, either union or intersection, to the tuples returned by the subqueries. For example, if we have $\underline{r}_1(\underline{a}_1, \underline{a}_2)$ with tuples $\{(1, 2), (3, 4)\}$ and $\underline{r}_2(\underline{a}_1, \underline{a}_2)$ with tuples $\{(1, 2), (5, 6)\}$ then $\underline{r}_1 \cup \underline{r}_2$ returns the tuples $\{(1, 2), (3, 4), (5, 6)\}$.

2.3 Variation Space in a Variational Database Framework

[have to revise] [define oplus in fig as syntactic sugar.]

To account for variability in a database we need to encode it. To encode variability we introduce a *feature space* \mathbf{F} as a closed set of boolean variables. A feature $f \in \mathbf{F}$ can be enabled (i.e., $f = \text{true}$) or disabled ($f = \text{false}$). Features capture the variation in a given variational scenario. For example, in the context of schema evolution, features can be generated from version numbers (e.g. features V_1 to V_5 and T_1 to T_5 in the motivating example, Table ??); for SPLs, the features can be adopted from the SPL feature set (e.g. the *edu* feature in our motivating example, Table ??); and for data integration, the features can be representatives of resources.

Features describe local points of variation in a database. Thus, enabling or disabling all the features produces a particular database *variant* where all variation has been removed. A *configuration* is a *total* function that maps every feature in the feature set to a boolean value. We represent a configuration by the set of enabled features. For example, in our motivating scenario, the configuration $\{V_2, T_3, \text{edu}\}$ represents a database variant where only features V_2 , T_3 , and *edu* are enabled. This database variant contains relation schemas in the yellow cells of Table ?. We refer to a variant with configuration c as variant c , e.g., variant $\{V_2, T_3, \text{edu}\}$ refers to the variant with configuration $\{V_2, T_3, \text{edu}\}$.

When describing variation points in the database, we need to refer to subsets of the configuration space. We achieve this by constructing propositional formulas of features. Thus, such a propositional formula defines a condition that holds for a subset of configurations and their corresponding variants. For example, the propositional formula $\neg \text{edu}$ represents all variants of our motivating example that do not have the education part of the schema, i.e., variant schemas of the left schema column.

We call a propositional formula of features a *feature expression* and define it formally in Figure ?. The evaluation function of feature expressions $\mathbb{E}[\![e]\!]_c : \mathbf{E} \rightarrow \mathbf{C} \rightarrow \mathbf{B}$ simply substitutes each feature f in the expression e with its boolean value set in the given configuration c and then simplifies the propositional formula to a boolean value. For example, assuming that $\mathbf{F} = \{f_1, f_2\}$ $\mathbb{E}[\![f_1 \vee f_2]\!]_{\{f_1\}} = \text{true} \vee \text{false} = \text{true}$, however, $\mathbb{E}[\![f_1 \vee f_2]\!]_{\{\}} = \text{false} \vee \text{false} = \text{false}$. Additionally, we define the binary *equivalence* (\equiv) relation and the unary *satisfiable* (*sat*) and *unsatisfiable* (*unsat*) relations over feature expressions in Figure ?.

To incorporate feature expressions into the database, we *annotate/tag* database elements (including attributes, relations, and tuples) with feature expressions. An *annotated element* x with feature expression e is denoted by x^e . The feature expression

Feature expression syntax:

$f \in \mathbf{F}$		<i>Feature Name</i>
$b \in \mathbf{B} := \text{true} \mid \text{false}$		<i>Boolean Value</i>
$e \in \mathbf{E} := b \mid f \mid \neg e \mid e \wedge e \mid e \vee e$		<i>Feature Expression</i>
$c \in \mathbf{C} : \mathbf{F} \rightarrow \mathbf{B}$		<i>Configuration</i>

Relations over feature expressions:

$$\begin{aligned}
e_1 \equiv e_2 & \text{ iff } \forall c \in \mathbf{C} : \mathbb{E}[e_1]_c = \mathbb{E}[e_2]_c \\
\text{sat}(e) & \text{ iff } \exists c \in \mathbf{C} : \mathbb{E}[e]_c = \text{true} \\
\text{unsat}(e) & \text{ iff } \forall c \in \mathbf{C} : \mathbb{E}[e]_c = \text{false}
\end{aligned}$$

Figure 2.2: Feature expression syntax and relations.

attached to an element is called a *presence condition* since it determines the condition (set of configurations) under which the element is present. $pc(x^e)$ returns the presence condition of the variational element x^e . For example, the annotated number $2^{f_1 \vee f_2}$ is present in variants $\{f_1\}$, $\{f_2\}$, $\{f_1, f_2\}$ but it does not exist in variant $\{\}$. Here, $pc(2^{f_1 \vee f_2}) = f_1 \vee f_2$.

No matter the context, features often have a relationship with each other that constrains configurations. For example, only one of the temporal features of $V_1 - V_5$ can be **true** for a given variant. This relationship is captured by a feature expression, called a *feature model* and denoted by m , which restricts the set of *valid configurations*: if configuration c violates the relationship then it follows from the $\mathbb{E}[\cdot]_c$ definition that $\mathbb{E}[m]_c = \text{false}$. For example, the restriction that at a given time only one of temporal features $V_1 - V_5$ can be enabled is represented by: $V_1 \oplus V_2 \oplus V_3 \oplus V_4 \oplus V_5$, where $f_1 \oplus f_2 \oplus \dots \oplus f_n$ is syntactic sugar for $(f_1 \wedge \neg f_2 \wedge \dots \wedge \neg f_n) \vee (\neg f_1 \wedge f_2 \wedge \dots \wedge \neg f_n) \vee (\neg f_1 \wedge \neg f_2 \wedge \dots \wedge f_n)$, i.e., features are mutually exclusive.

2.4 Variational Set

[havve to revise this]

A *variational set* (*v-set*) $X = \{x_1^{e_1}, \dots, x_n^{e_n}\}$ is a set of annotated elements $[?, ?, ?]$. We typically omit the presence condition **true** in a variational set, e.g., $4^{\text{true}} = 4$.

Semantics of feature expressions:

$$\begin{aligned}
\mathbb{E}[\cdot] : \mathbf{E} &\rightarrow \mathbf{C} \rightarrow \mathbf{B} \\
\mathbb{E}[b]_c &= b \\
\mathbb{E}[f]_c &= c \cdot f \\
\mathbb{E}[\neg f]_c &= \neg \mathbb{E}[f]_c \\
\mathbb{E}[e_1 \wedge e_2]_c &= \mathbb{E}[e_1]_c \wedge \mathbb{E}[e_2]_c \\
\mathbb{E}[e_1 \vee e_2]_c &= \mathbb{E}[e_1]_c \vee \mathbb{E}[e_2]_c
\end{aligned}$$

Figure 2.3: Feature expression evaluation and functions.

Conceptually, a *variational set* represents many different plain sets simultaneously. Thus, a plain set $\underline{X} \in \underline{\mathcal{X}}$ is a *variant* of the variational set $X \in \mathcal{X}$ and given variant \underline{X} 's configuration $c \in \mathbf{C}$ it can be generated by evaluating the presence condition of each element with c and including the element if the said evaluation results in **true** and excluding it otherwise. This is done by the *v-set configuration* function $\mathbb{X}[X]_c : \mathcal{X} \rightarrow \mathbf{C} \rightarrow \underline{\mathcal{X}}$. For example, assume we have the feature space $\mathbf{F} = \{f_1, f_2\}$ and the v-set $X_1 = \{2^{f_1}, 3^{f_2}, 4\}$. X_1 represents four plain sets: $\{2, 3, 4\} = \mathbb{X}[X_1]_{\{f_1, f_2\}}$, i.e., the variant $\{2, 3, 4\}$ is generated from the v-set X_1 under the configuration $f_1 = \mathbf{true}, f_2 = \mathbf{true}$; $\{2, 4\} = \mathbb{X}[X_1]_{\{f_1\}}$; $\{3, 4\} = \mathbb{X}[X_1]_{\{f_2\}}$; and $\{4\} = \mathbb{X}[X_1]_{\{\}}$. Following database notational conventions we drop the brackets of a variational set when used in database schema definitions and queries.

A variational set itself can also be annotated with a feature expression. $X^e = \{x_1^{e_1}, \dots, x_n^{e_n}\}^e$ is an *annotated v-set*. Annotating a v-set with the feature expression e restricts the condition under which its elements are present, i.e., it forces elements' presence conditions to be more specific. The *normalization* operation $\downarrow(X^e)$ applies this restriction: $\downarrow(X^e) = \{x_i^{e_i \wedge e} \mid x_i^{e_i} \in X^e, \text{sat}(e_i \wedge e)\}$. Note that the *normalization* operation also removes the elements with unsatisfiable presence conditions and it can also be applied to a v-set X since $X^{\mathbf{true}} = X$. For example, the annotated v-set $X_1 = \{2^{f_1}, 3^{\neg f_2}, 4, 5^{f_3}\}^{f_1 \wedge f_2}$ indicates that all the elements of the set can only exist when both f_1 and f_2 are enabled. Thus, normalizing the v-set X_1 results in $\{2^{f_1 \wedge f_2}, 4^{f_1 \wedge f_2}, 5^{f_1 \wedge f_2 \wedge f_3}\}$. The element 3 is dropped since $\neg \text{sat}(pc(3, X_1))$, where

$pc(3, X_1) = \neg f_2 \wedge (f_1 \wedge f_2)$. Note that we use the function $pc(x, X^e)$ to return the presence condition of a unique variational element within a bigger variational structure.

We provide some operations over v-sets used mainly in the Section ?? and defined formally in Appendix ?. Intuitively, these operations should behave such that configuring the result of applying a variational set operation should be equivalent to applying the plain set operation on the configured input v-sets for all valid configurations. Thus, for all possible operations, \odot , defined on v-sets the property $(P_1) : \forall c \in \mathbf{C}. \mathbb{X}[\downarrow(X_1) \odot \downarrow(X_2)]_c = \mathbb{X}[\downarrow(X_1)]_c \odot \mathbb{X}[\downarrow(X_2)]_c$ must hold, where \odot is the counterpart operation on plain sets.

Definition 2.4.1 (V-set union) *The union of two v-sets is the union of their elements with the disjunction of presence conditions if an element exists in both v-sets: $X_1 \cup X_2 = \{x^{e_1} \mid x^{e_1} \in \downarrow(X_1), \exists e_2. x^{e_2} \notin \downarrow(X_2)\} \cup \{x^{e_2} \mid x^{e_2} \in \downarrow(X_2), \exists e_1. x^{e_1} \notin \downarrow(X_1)\} \cup \{x^{e_1 \vee e_2} \mid x^{e_1} \in \downarrow(X_1), x^{e_2} \in \downarrow(X_2)\}$. For example, the union of two normalized v-sets $\{2, 3^{e_1}, 4^{e_1}\}$ and $\{3^{e_2}, 4^{\neg e_1}\}$ is $\{2, 3^{e_1 \vee e_2}, 4\}$.*

Definition 2.4.2 (V-set intersection) *The intersection of two v-sets is a v-set of their shared elements annotated with the conjunction of their presence conditions, i.e., $X_1 \cap X_2 = \{x^{e_1 \wedge e_2} \mid x^{e_1} \in X_1, x^{e_2} \in X_2, \text{sat}(e_1 \wedge e_2)\}$.*

Definition 2.4.3 (V-set cross product) *The cross product of two v-sets is a pair of every two elements of them annotated with the conjunction of their presence conditions. $X_1 \times X_2 = \{(x_1, x_2)^{e_1 \wedge e_2} \mid x_1^{e_1} \in X_1, x_2^{e_2} \in X_2\}$*

Definition 2.4.4 (V-set equivalence) *Two v-sets are equivalent, denoted by $X_1 \equiv X_2$, iff $\forall x^e \in (\downarrow(X_1) \cup \downarrow(X_2)), x^{e_1} \in \downarrow(X_1), x^{e_2} \in \downarrow(X_2). e_1 \equiv e_2, e \equiv e_1$, i.e., they both cover the same set of elements and the presence conditions of elements from the two normalized v-sets are equivalent.*

2.4.1 Variational Set Configuration

[vset configuration.]

2.5 The Formula Choice Calculus

[formula choice calculus]

The choice calculus $[?, ?]$ is a metalanguage for describing variation in programs and its elements such as data structures $[?, ?]$. In the choice calculus, variation is represented in-place as choices between alternative subexpressions. For example, the variational expression $e = A\langle 1, 2 \rangle + B\langle 3, 4 \rangle + A\langle 5, 6 \rangle$ contains three choices. Each choice has an associated *dimension*, which is used to synchronize the choice with other choices in different parts of the expression. For example, expression e contains two dimensions, A and B , and the two choices in dimension A are synchronized. Therefore, the variational expression e represents four different plain expressions, depending on whether the left or right alternatives are selected from each dimension. Assuming that dimensions may be set to boolean values where `true` indicates the left alternative and `false` indicates the right alternative, we have: (1) $1 + 3 + 5$, when A and B are `true`, (2) $1 + 4 + 5$, when A is `true` and B is `false`, (3) $2 + 3 + 6$, when A is `false` and B is `true`, and (4) $2 + 4 + 6$, when A and B are `false`.

The formula choice calculus $[?]$ extends the choice calculus by allowing dimensions to be propositional formulas. For example, the variational expression $e' = A \vee B\langle 1, 2 \rangle$ represents two plain expressions: (1) 1 , when $A \vee B$ evaluates to `true` and (2) 1 , when $A \vee B$ evaluates to `false`.

A formula e can also be used to *annotate/tag* an element x which is denoted by x^e . For example, the variational expression $1^{A \vee B}$ states that the expression 1 is present when $A \vee B$ evaluates to `true` and otherwise it is absent.

Chapter 3: The Variational Database Framework

[needs. must have configuration.]

3.1 Variational Schema

[vsch]

3.1.1 Variational Schema Configuration

[vsch configuration.]

3.2 Variational Table

[vtab]

3.2.1 Variational Table Configuration

[vtab configuration]

3.3 Variational Database

[vdb]

3.3.1 Variational Database Configuration

[vdb configuration]

3.4 Properties of a Variational Database Framework

[well-formed vdb properties.context-specific properties.]

[show that they hold for vdb.]

Chapter 4: The Variational Query Language

[vql]

4.1 Variational Relational Algebra

[vra]

4.1.1 VRA Configuration

[vra configuration]

4.1.2 VRA Semantics

[vra semantics]

4.1.3 VRA Type System

[type sys]

4.1.4 VRA Variation-Minimization Rules

[rules]

4.2 Variational Query Language Properties

[prop. show for vra.]

Chapter 5: Variational Database Usecases

5.1 Variation in Space

[enron email usecase]

5.2 Variation in Time

[employee evolution usecase]

Chapter 6: Variational Database Management System (VDBMS)

[vdbms]

6.1 Implemented Approaches

[apps]

6.2 Experiments

[exp.]

Chapter 7: Related Work

[related work! have to work on this!]

7.1 Instances of Variation in Databases

[schema evolution. database versioning. data integration. data provenance.]

7.2 Instances of Database Variation Resulted from Software Development

[SPL. data model. query.]

7.3 Variational Research

[blah]

Chapter 8: Conclusion

[conclusion]

Bibliography

- [1] Serge Abiteboul, Richard Hull, and Victor Vianu. *Foundations of Databases: The Logical Level*. Addison-Wesley, 1994.
- [2] Lamia Abo Zaid and Olga De Troyer. Towards modeling data variability in software product lines. In Terry Halpin, Selmin Nurcan, John Krogstie, Pnina Soffer, Erik Proper, Rainer Schmidt, and Ilia Bider, editors, *Enterprise, Business-Process and Information Systems Modeling*, pages 453–467, Berlin, Heidelberg, 2011. Springer.
- [3] Gad Ariav. Temporally oriented data definitions: Managing schema evolution in temporally oriented databases. *Data & Knowledge Engineering*, 6(6):451 – 467, 1991.
- [4] Parisa Ataei, Qiaoran Li, and Eric Walkingshaw. Should variation be encoded explicitly in databases? In *15th International Working Conference on Variability Modelling of Software-Intensive Systems*, VaMoS’21, New York, NY, USA, 2021. Association for Computing Machinery.
- [5] Parisa Ataei, Arash Termehchy, and Eric Walkingshaw. Variational Databases. In *Int. Symp. on Database Programming Languages (DBPL)*, pages 11:1–11:4. ACM, 2017.
- [6] Anant P. Bhardwaj, Souvik Bhattacharjee, Amit Chavan, Amol Deshpande, Aaron J. Elmore, Samuel Madden, and Aditya G. Parameswaran. Datahub: Collaborative data science & dataset version management at scale. In *CIDR 2015, Seventh Biennial Conference on Innovative Data Systems Research, Asilomar, CA, USA, January 4-7, 2015, Online Proceedings*. www.cidrdb.org, 2015.
- [7] Souvik Bhattacharjee, Amit Chavan, Silu Huang, Amol Deshpande, and Aditya Parameswaran. Principles of dataset versioning: Exploring the recreation/storage tradeoff. *Proc. VLDB Endow.*, 8(12):1346–1357, August 2015.
- [8] Cristina De Castro, Fabio Grandi, and Maria Rita Scalas. Schema versioning for multitemporal relational databases. *Information Systems*, 22(5):249 – 290, 1997.
- [9] Badrish Chandramouli, Johannes Gehrke, Jonathan Goldstein, Donald Kossmann, Justin J. Levandoski, Renato Marroquin, and Wenlei Xie. READY: completeness

- is in the eye of the beholder. In *CIDR 2017, 8th Biennial Conference on Innovative Data Systems Research, Chaminade, CA, USA, January 8-11, 2017, Online Proceedings*. www.cidrdb.org, 2017.
- [10] Paul Clements and Linda Northrop. *Software Product Lines: Practices and Patterns*. Addison-Wesley, Boston, MA, 2001.
 - [11] AnHai Doan, Alon Halevy, and Zachary Ives. *Principles of Data Integration*. Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, 1st edition, 2012.
 - [12] Martin Erwig and Eric Walkingshaw. The Choice Calculus: A Representation for Software Variation. *ACM Trans. on Software Engineering and Methodology (TOSEM)*, 21(1):6:1–6:27, 2011.
 - [13] Martin Erwig, Eric Walkingshaw, and Sheng Chen. An Abstract Representation of Variational Graphs. In *Int. Work. on Feature-Oriented Software Development (FOSD)*, pages 25–32. ACM, 2013.
 - [14] Mina Farid, Alexandra Roatis, Ihab F. Ilyas, Hella-Franziska Hoffmann, and Xu Chu. Clams: Bringing quality to data lakes. In *Proceedings of the 2016 International Conference on Management of Data*, SIGMOD ’16, page 2089–2092, New York, NY, USA, 2016. Association for Computing Machinery.
 - [15] Kai Herrmann, Jan Reimann, Hannes Voigt, Birgit Demuth, Stefan Fromm, Robert Stelzmann, and Wolfgang Lehner. Database evolution for software product lines. In *DATA*, 2015.
 - [16] Silu Huang, Liqi Xu, Jialin Liu, Aaron J. Elmore, and Aditya Parameswaran. Orpheusdb: Bolt-on versioning for relational databases. *Proc. VLDB Endow.*, 10(10):1130–1141, June 2017.
 - [17] Spencer Hubbard and Eric Walkingshaw. Formula Choice Calculus. In *Int. Work. on Feature-Oriented Software Development (FOSD)*, pages 49–57. ACM, 2016.
 - [18] Edwin McKenzie and Richard Thomas Snodgrass. Schema evolution and the relational algebra. *Inf. Syst.*, 15(2):207–232, May 1990.
 - [19] Hyun J. Moon, Carlo A. Curino, Alin Deutsch, Chien-Yi Hou, and Carlo Zaniolo. Managing and querying transaction-time databases under schema evolution. *Proc. VLDB Endow.*, 1(1):882–895, August 2008.
 - [20] Martin Schäler, Thomas Leich, Marko Rosenmüller, and Gunter Saake. Building information system variants with tailored database schemas using features. In *Jolita*

- Ralyté, Xavier Franch, Sjaak Brinkkemper, and Stanislaw Wrycza, editors, *Advanced Information Systems Engineering*, pages 597–612, Berlin, Heidelberg, 2012. Springer.
- [21] Norbert Siegmund, Christian Kästner, Marko Rosenmüller, Florian Heidenreich, Sven Apel, and Gunter Saake. Bridging the Gap Between Variability in Client Application and Database Schema. In *13. GI-Fachtagung Datenbanksysteme für Business, Technologie und Web (BTW)*, pages 297–306. Gesellschaft für Informatik (GI), 2009.
 - [22] Richard Thomas Snodgrass. *The TSQL2 Temporal Query Language*. Kluwer Academic Publishers, USA, 1995.
 - [23] Micheal Stonebraker, Dong Deng, and Micheal L. Brodie. Database decay and how to avoid it. In *Big Data (Big Data), 2016 IEEE International Conference*. IEEE, 2016.
 - [24] U. Störl, D. Müller, A. Tekleab, S. Tolale, J. Stenzel, M. Klettke, and S. Scherzinger. Curating variational data in application development. In *2018 IEEE 34th International Conference on Data Engineering (ICDE)*, pages 1605–1608, 2018.
 - [25] Thomas Thüm, Leopoldo Teixeira, Klaus Schmid, Eric Walkingshaw, Mukelabai Mukelabai, Mahsa Varshosaz, Goetz Botterweck, Ina Schaefer, and Timo Kehler. Toward Efficient Analysis of Variation in Time and Space. In *Int. Work. on Variability and Evolution of Software Intensive Systems (VariVolution)*, 2019.
 - [26] Eric Walkingshaw. *The Choice Calculus: A Formal Language of Variation*. PhD thesis, Oregon State University, 2013. <http://hdl.handle.net/1957/40652>.
 - [27] Eric Walkingshaw, Christian Kästner, Martin Erwig, Sven Apel, and Eric Bodden. Variational Data Structures: Exploring Trade-Offs in Computing with Variability. In *ACM SIGPLAN Symp. on New Ideas in Programming and Reflections on Software (Onward!)*, pages 213–226, 2014.

APPENDICES

