

Advancing Fact Attribution for Query Answering: Aggregate Queries and Novel Algorithms

Omer Abramovich
Tel Aviv University
Tel Aviv, Israel

Daniel Deutch
Tel Aviv University
Tel Aviv, Israel

Nave Frost
eBay Research
Netanya, Israel

Ahmet Kara
OTH Regensburg
Regensburg, Germany

Dan Olteanu
University of Zurich
Zurich, Switzerland

ABSTRACT

In this paper, we introduce a novel approach to computing the contribution of input tuples to the answers of queries, quantified using the Banzhaf and Shapley values. In contrast to prior algorithmic work focusing on Select-Project-Join-Union queries, ours is the first practical approach for queries with aggregates.

Our approach relies on two novel optimizations that are essential for its practicality, and also significantly improve the performance even for queries without aggregates. The first optimization exploits the observation that many input tuples have the same contributions to the query answer, so it is enough to compute the contribution of one of them. The second optimization uses the gradient of the query lineage to compute the contributions of all tuples in the time needed for one of them. Experiments show that our approach achieves up to 3 orders of magnitude runtime improvements over the state-of-the-art for queries without aggregates, and that it is indeed practical for aggregate queries.

PVLDB Reference Format:

Omer Abramovich, Daniel Deutch, Nave Frost, Ahmet Kara, and Dan Olteanu. Advancing Fact Attribution for Query Answering: Aggregate Queries and Novel Algorithms. PVLDB, 14(1): XXX-XXX, 2020.
doi:XX.XX/XXX.XX

PVLDB Artifact Availability:

The source code, data, and/or other artifacts have been made available at <https://github.com/Omer-Abramovich/LExaBan-LExaShap>.

1 INTRODUCTION

Recent years have witnessed a surge of development that uses game theoretic measures to quantify the contribution of each database tuple to the query answer, also referred to as *attribution*. The Banzhaf [4, 27] and Shapley [31] values are two prominent examples of such measures. They originate in cooperative game theory, where they are used to quantify the contribution of a player by summing up their marginal contributions over the subsets (or permutations, for Shapley) of players. They are used in applications across several domains. For instance, in machine learning, they are

used for feature importance and data valuation, feature selection, and explanations [29, 33]. Prior works used Banzhaf and Shapley values for data attribution in query answering: the contribution of a database tuple to the query answer is quantified by defining a game where the players are database tuples and each tuple in the query result is treated as the objective function [2, 11, 19]. Banzhaf and Shapley values are two different measurements, yet they are closely related, and in particular previous work [2] has shown in the context of query answering they tend to agree on the relative order of fact contribution.

Example 1.1. Figure 1 depicts an IMDB-like database with information on some movies directed by Tarantino.

Query Q_1 asks whether Tarantino was nominated for the academy awards for "best director". The query answer is true, since Tarantino was nominated for 3 movies. We would like to understand which leading actor contributed the most to the query answer which in this case corresponds to Tarantino's nominated movies. To measure such contribution, we can use importance scores such as the Banzhaf or Shapley values. In our example, it turns out that the actor with the largest (Banzhaf and Shapley) score is Brad Pitt, having played a leading role in 2 of these 3 movies.

Query Q_2 asks for the maximum revenue of a movie directed by Tarantino. Again, we would like to understand which leading actor contributed the most to the query answer, i.e., to Tarantino's high-revenue movies. Using Banzhaf or Shapley values, we find that Samuel L. Jackson contributed the most, followed by Leonardo DiCaprio, Brad Pitt, Uma Thurman, and Christoph Waltz.

Multiple previous works have demonstrated the usefulness of Banzhaf and Shapley values [2, 6, 7, 16–19, 21, 28]. In this paper, we study the problem of *computing the contribution of database tuples to the answers to aggregate queries, quantified using the Banzhaf and Shapley values*. The computational problem is known to be highly intractable in general: already for Boolean non-hierarchical queries (which form a strict subclass of acyclic queries), the exact computation of Shapley and Banzhaf values is #P-hard [19]. Furthermore, ranking tuples by their Banzhaf or Shapley values remains intractable under standard theoretic assumptions [2]. While previous work has proposed multiple optimization strategies for computing Banzhaf and Shapley values in the context of SPJU queries, the algorithms' execution time is far from negligible, especially when considering computation of Banzhaf/Shapley scores for all input fact-output tuple pairs. *Furthermore, all practical algorithms and prototype implementations in this context are limited in their*

This work is licensed under the Creative Commons BY-NC-ND 4.0 International License. Visit <https://creativecommons.org/licenses/by-nc-nd/4.0/> to view a copy of this license. For any use beyond those covered by this license, obtain permission by emailing info@vldb.org. Copyright is held by the owner/author(s). Publication rights licensed to the VLDB Endowment.

Proceedings of the VLDB Endowment, Vol. 14, No. 1 ISSN 2150-8097.
doi:XX.XX/XXX.XX

DirectingAwards (endo)

	Award
d_1	Academy
d_2	BAFTA

Movies (endo)

	Title	Director	Gross revenue(M\$)
m_1	Kill Bill: Vol. 1	Tarantino	176
m_2	Inglourious Basterds	Tarantino	322
m_3	Once Upon a Time in Hollywood	Tarantino	377

Cast (endo)

	Name
a_1	Brad Pitt
a_2	Leonardo DiCaprio
a_3	Zoë Bell
a_4	Uma Thurman

MovieAwards (exo)

Movie ID	Award
Inglourious Basterds	Academy
Inglourious Basterds	BAFTA
Once Upon a Time in Hollywood	Academy
Once Upon a Time in Hollywood	BAFTA

Query Q_1

$$Q() = \text{Movies}(x, \text{Tarantino}, r),$$

$$\text{Cast}(y),$$

$$\text{DirectingAwards}(z),$$

$$\text{MovieCast}(y, x),$$

$$\text{MovieAwards}(x, z)$$
Query Q_2

$$\langle \text{MAX}, r, \rangle$$

$$Q(x, y, r) =$$

$$\text{Movies}(x, \text{Tarantino}, r),$$

$$\text{Cast}(y),$$

$$\text{MovieCast}(y, x)\rangle$$
MovieCast (exo)

Movie ID	Actor ID
Kill Bill: Vol. 1	Uma Thurman
Inglourious Basterds	Brad Pitt
Once Upon a Time in Hollywood	Brad Pitt
Inglourious Basterds	Zoë Bell
Once Upon a Time in Hollywood	Zoë Bell
Once Upon a Time in Hollywood	Leonardo DiCaprio

Figure 1: Running example: an IMDB-like database and two queries

applicability, in that they have focused on SQL queries with select, project, join, and union.

In this paper, we go beyond prior work by adding support to aggregates (COUNT, SUM, MIN, MAX) in the select clause of the query. Aggregation is key to practical OLAP and decision support workloads; all 22 TPC-H queries involve aggregation.

All prior practical approaches to fact attribution, which were designed for queries without aggregates, do not work for queries with aggregates. Such approaches compile the lineage (also called provenance) of the query [8, 14, 25] into tractable circuits that support the efficient computation of the Banzhaf and Shapley values for the input tuples. Yet the lineage formalism necessarily differs between the two query classes: Whereas the lineage for the former is a *Boolean semiring* expression, for the latter it is a *semimodule* expression, which pairs an *aggregate monoid* expression with a Boolean semiring expression. *The first main contribution of this paper is an extension of the compilation approach for Banzhaf- and Shapley-based fact attribution from queries without aggregates to queries with aggregates, so from semiring to semimodule expressions.*

Yet, the computation of Banzhaf/Shapley values for semimodule expressions remains expensive for most benchmarks used in the literature [2, 11]. *The second main contribution of this paper is then an interplay of two novel techniques that we incorporate in our semimodule compilation algorithm: variable lifting and Banzhaf/Shapley gradient.* These techniques significantly reduce the computation time for fact attribution for *queries with and even without aggregates.*

In more detail, our contributions are as follows. We introduce LEXABAN and LEXASHAP, two novel algorithms to compute Banzhaf and respectively Shapley values for SPJUA (Select-Project-Join-Union-Aggregate) queries. We conduct experiments with these algorithms and the state-of-the-art algorithms EXABAN and EXASHAP for fact attribution and SPJU queries using more than a million instances from 3 databases [2]. For queries without aggregates, LEXABAN achieves 2-3 orders of magnitude improvement over EXABAN and LEXASHAP achieves more than 2 orders of magnitude improvement over EXASHAP. For queries with aggregates, LEXABAN and LEXASHAP are, to our knowledge, the first practical algorithms for Banzhaf and Shapley value computation. The performance gains of LEXABAN and LEXASHAP over the state-of-the-art are primarily due to two novel techniques.

Variable Lifting. A major challenge faced by state-of-the-art approaches to fact attribution in query answering is the need to compile query lineage into forms that admit tractable Banzhaf/Shapley value computation. A prime example of such tractable form is

the decomposition tree (d-trees) [12]: This is a Boolean circuit (1) that has gates expressing independence or mutual exclusiveness of their children and (2) on which the Banzhaf/Shapley values can be computed in one bottom-up pass. The compilation can be time consuming and often yields large d-trees.

We observe that some variables in the lineage are interchangeable. For instance, in the Boolean formula $(x_1 \wedge y_1 \wedge z_1) \vee (x_2 \wedge y_1 \wedge z_1)$ the pairs of variables $\{x_1, x_2\}$ and $\{y_1, z_1\}$ are interchangeable. We extend the d-trees with a new gate type, which replace, or *lift*, conjunctions or disjunctions of variables that behave the same in the lineage with a fresh variable. Such lift gates preserve the equivalence to the original formula. Variable lifting can drastically reduce the compilation time and lead to much smaller d-trees. In our experiments, it leads to more than 2 OOM (orders-of-magnitude) speedup and more than 1 OOM smaller d-trees.

Gradient-based Banzhaf/Shapley value computation. A further major challenge is that we need to compute the Banzhaf/Shapley values for all input facts. Computing such values separately for each fact is computationally expensive. We exhibit a novel relationship between the Banzhaf/Shapley values and the gradients of a function defined by the query lineage. This allows us to devise back-propagation algorithms that compute these values efficiently over the d-tree. Our gradient technique allows to compute the Banzhaf/Shapley values for all variables of a given query lineage in the same time needed to compute them for one variable. In our experiments, this gradient technique speeds up the computation by a factor that grows with the size of the instance and up to more than 2 OOM speedup for large instances. This enables the computation of Banzhaf values for instances larger than previously possible.

We first introduce these two techniques for SPJU queries in Sec. 3 and then extend them to SPJUA queries in Sec. 4. In our experiments (Sec. 5), we show that the interplay of both techniques can lead to 3 OOM speedup.

Related Work. Banzhaf and Shapley values, originally introduced in the context of game theory to quantify the contribution of individual players to the value of a cooperative game [4, 31] were used in query answering to quantify the contribution of individual database facts to query answers [19]. Previous work has analyzed the complexity of computing these measures [19, 28] and proposed efficient algorithms to compute them based on knowledge compilation algorithms [2, 11, 17]. Since the problem is computationally intractable in general, multiple heuristics and approximation schemes were proposed. Our algorithms are far more efficient than state of the art algorithms, pushing the boundaries of tractable instances.

Many other notions have been used to quantify fact contribution in query answering, including causality [25], responsibility [23] or counterfactuals [24]. Banzhaf and Shapley values are unique in that they are grounded in game theory where their properties have been extensively studied, and previous work has shown their usefulness for query answering. Banzhaf and Shapley values are prevalent in many other fields, notably interpreting model predictions and feature selection in Machine Learning [13, 32]. For instance, the SHAP score [20] leverages Shapley values to explain model predictions. These methods differ from the query answering scenario in that "missing players" are modeled according to their expected values.

A full version of the paper, including proofs for all formal results, as well as the implementation of all algorithms, are available in [3].

2 PRELIMINARIES

We use \mathbb{N} to denote the set of natural numbers including 0. For $n \in \mathbb{N}$, we define $[n] = \{1, \dots, n\}$. In case $n = 0$, we have $[n] = \emptyset$.

Boolean Formulas. A *Boolean formula* (or simply formula) over a set X of Boolean variables is either a constant 0 or 1, a variable $x \in X$, a negation $\neg\varphi$, or a conjunction/disjunction $\varphi_1 \wedge \varphi_2$ or $\varphi_1 \vee \varphi_2$ of two formulas φ_1, φ_2 . We denote by $\text{Bool}(X)$ the set of Boolean formulas over X . A *literal* is a constant, a variable, or its negation. The set of variables in φ is denoted $\text{vars}(\varphi)$, and φ is *read-once* if each variable appears at most once. Substituting $x \in X$ with $b \in \{0, 1\}$ in φ is written as $\varphi[x := b]$. A *valuation* $\theta : \text{vars}(\varphi) \rightarrow \{0, 1\}$ for φ maps each variable in φ to 0 or 1, and we identify θ with $\{x \in \text{vars}(\varphi) \mid \theta(x) = 1\}$, where $|\theta|$ is the number of variables mapped to 1. The Boolean value of φ under θ is $\varphi[\theta]$, and θ is a *model* of φ if $\varphi[\theta] = 1$. The *model set* $\text{models}(\varphi)$ consists of all models of φ , with *model count* $\#\varphi = |\text{models}(\varphi)|$ and *k-model count* $\#_k\varphi = |\text{models}(\varphi) \cap \binom{X}{k}|$, where $\binom{X}{k}$ is the set of k -element subsets of X . Two formulas φ_1, φ_2 are *independent* if $\text{vars}(\varphi_1) \cap \text{vars}(\varphi_2) = \emptyset$ and *exclusive* if $\text{models}(\varphi_1) \cap \text{models}(\varphi_2) = \emptyset$. The set $\text{PosBool}(X)$ of *positive formulas* (formulas without negation) forms a commutative semiring under $\vee, \wedge, 0, 1$ [15].

Example 2.1. Consider the formula $\varphi = x \vee y$. The models of φ are $\{x\}$, $\{y\}$, and $\{x, y\}$. Hence, $\#\varphi = 3$, $\#_0\varphi = 0$, $\#_1\varphi = 2$, and $\#_2\varphi = 1$.

Databases. We follow standard relational database notions [1]. We assume a countably infinite set Const of constants used as database values. For a tuple $t \in \text{Const}^k$, $k \in \mathbb{N}$ and $i \in [k]$, we denote its i -th value by $t[i]$. A database consists of *facts* $R(t)$, where R is a relation name and t is a tuple matching R 's arity. Following prior work, each database D is partitioned into *endogenous* facts D_n and *exogenous* facts D_x [19].

Conjunctive Queries. A *conjunctive query* (CQ) has the form $Q(X) = R_1(X_1), \dots, R_m(X_m)$, where each R_i is a relation name; each X_i is a tuple of variables and constants matching the arity of R_i . The *head variables* X form a subset of the variables in X_1, \dots, X_m . A CQ is *k-ary* if $|X| = k$ and *Boolean* if $k = 0$. Each $R_i(X_i)$ is called an *atom*. A *union of conjunctive queries* (UCQ) consists of a set of k -ary CQs for some $k \in \mathbb{N}$.

Example 2.2. Query Q_1 in Figure 1 is a CQ asking whether movies directed by Tarantino were nominated for directing awards.

Aggregate Queries. We define *aggregate queries* following prior work [19]. Given a database D , an aggregate query is a triple $\langle \alpha, \gamma, Q \rangle$, where Q is a k -ary UCQ for some $k \in \mathbb{N}$, $\gamma : \text{Const}^k \rightarrow \mathbb{R}$ is a function that maps each k -ary value tuple to a numeric value and $\alpha : \mathbb{R}^* \rightarrow \mathbb{R}$ is a function that maps a bag of numeric values to one value. Given a database D , let $Q(D) = \{t_1, \dots, t_k\}$. The result of $\langle \alpha, \gamma, Q \rangle$ over D is defined as $\langle \alpha, \gamma, Q \rangle(D) := \alpha(\gamma(t_i)_{i \in [k]})$. If γ returns the i -th value of the tuple t , we just write i instead of α . We consider SUM and MAX aggregate queries of the following form:

$$\langle \text{SUM}, \gamma, Q \rangle(D) \stackrel{\text{def}}{=} \sum_{\vec{c} \in Q(D)} \gamma(\vec{c})$$

$$\langle \text{MAX}, \gamma, Q \rangle(D) \stackrel{\text{def}}{=} \begin{cases} \max \{ \gamma(\vec{c}) \mid \vec{c} \in Q(D) \} & \text{if } Q(D) \neq \emptyset \\ 0 & \text{if } Q(D) = \emptyset \end{cases}$$

The query $\langle \text{MIN}, \gamma, Q \rangle$ is defined analogously to $\langle \text{MAX}, \gamma, Q \rangle$.

We define $\langle \text{COUNT}, Q \rangle \stackrel{\text{def}}{=} \langle \text{SUM}, 1, Q \rangle$, where 1 maps each output tuple to 1, yielding the number of output tuples of Q . If the aggregate is clear from the context, we refer to an aggregate query by Q .

Remark 2.3. The choice for the MAX and MIN queries to evaluate to 0 if $Q(D) = \emptyset$ is not arbitrary. The alternatives, using ∞ for MIN and $-\infty$ for MAX, are not useful in the context of attribution values, in which we sum up over different valuations of the query.

Example 2.4. Query Q_2 of Figure 1 is an aggregate query that asks for the maximal revenue of a movie directed by Tarantino.

Monoids and Semimodules. The development in this paper uses the algebraic structures of numeric monoids and semimodules constructed using such monoids and the $\text{PosBool}(X)$ semiring. We denote by $\overline{\mathbb{R}}$ the set of real numbers \mathbb{R} including $\{\infty, -\infty\}$.

Definition 2.5. A numeric monoid (or simply, a monoid) $(\overline{\mathbb{R}}, +_M, 0_M)$ consists of a binary operation $+_M : \overline{\mathbb{R}} \times \overline{\mathbb{R}} \rightarrow \overline{\mathbb{R}}$, and a neutral element $0_M \in \overline{\mathbb{R}}$ s.t. for all $m_1, m_2, m_3 \in \overline{\mathbb{R}}$, the following holds:

$$\begin{aligned} (m_1 +_M m_2) +_M m_3 &= m_1 +_M (m_2 +_M m_3) \\ 0 +_M m_1 &= m_1 +_M 0 = m_1 \end{aligned}$$

A monoid is *commutative* if $m_1 +_M m_2 = m_2 +_M m_1$ for all $m_1, m_2 \in M$.

Monoids naturally model many aggregate functions. For instance, the SQL aggregate function MAX, which returns the maximum value in a relation column is modeled by the monoid $(\max, -\infty)$.

Next, we introduce PosBool -semimodules, which combine monoid values with positive Boolean formulas to represent numerical values conditioned by Boolean formulas.

Definition 2.6. Given a monoid $M = (\overline{\mathbb{R}}, +_M, 0_M)$ and the semiring $(\text{PosBool}(X), \vee, \wedge, 0, 1)$, the triple $(\text{PosBool}(X), M, \otimes)$ is a PosBool -semimodule where $\otimes : \text{PosBool}(X) \times \overline{\mathbb{R}} \rightarrow \overline{\mathbb{R}}$, such that for all $\varphi_1, \varphi_2 \in \text{PosBool}(X)$ and $m_1, m_2 \in \overline{\mathbb{R}}$, the following axioms hold:

- (1) $\varphi_1 \otimes (m_1 +_M m_2) = \varphi_1 \otimes m_1 +_M \varphi_1 \otimes m_2$
- (2) $(\varphi_1 \vee \varphi_2) \otimes m_1 = \varphi_1 \otimes m_1 +_M \varphi_2 \otimes m_1$
- (3) $(\varphi_1 \wedge \varphi_2) \otimes m_1 = \varphi_1 \otimes (\varphi_2 \otimes m_1)$
- (4) $\varphi_1 \otimes 0_M = 0 \otimes m_1 = 0_M$
- (5) $1 \otimes m_1 = m_1$

A *semimodule expression* is either an element from M or of one of the forms $\Phi_1 +_M \Phi_2$ and $\varphi \otimes \Phi_1$ for semimodule expressions Φ_1 and Φ_2 and a Boolean formula $\varphi \in \text{PosBool}(X)$. It easily follows from the definition of semimodules that every semimodule expression

can be equivalently written as $\varphi_1 \otimes m_1 +_M \dots +_M \varphi_k \otimes m_k$, where each φ_i is a Boolean formula and each m_i is an element from M . We abbreviate the latter expression by $\sum_M \varphi_i \otimes m_i$.

Consider a semimodule expression $\Phi = \sum_M \varphi_i \otimes m_i$. The *Boolean part* of Φ is $\bigvee \varphi_i$. Given a valuation θ over X , we denote $S_\theta = \{i | \varphi_i[\theta] = 1\}$. By $\Phi[\theta]$, we denote the monoid value computed by $\sum_M [\varphi_j[\theta] \otimes m_j]_{j \in S_\theta}$. Note that this is equal to $\sum_M m_j$ where the sum ranges over all j for which φ_j evaluates to 1 over θ .

Remark 2.7. Note that the choice to evaluate the semimodule expression over formulas φ_i s.t. $\varphi_i[\theta] = 1$ instead of evaluating it over all formulas is not standard, and is done to accommodate for the definition of the MAX and MIN queries as defined for attribution.

For a semimodule expression Φ , we denote by $\#^P \Phi$ the number of valuations θ such that $\Phi[\theta] = p$ and by $\#_k^P \Phi$ the number of such valuations of size k . We define the variables $\text{vars}(\Phi)$ of Φ to be the variables $\text{vars}(\varphi)$ of φ and the model count $\#\Phi$ of Φ to be the model count $\#\varphi$ of φ , where φ is the Boolean part of Φ . We say that two semimodule expressions T_1 and T_2 are independent (resp. exclusive) if their Boolean parts are independent (resp. exclusive).

Query Grounding. A *grounding* of a CQ Q w.r.t. a database D is a mapping G of the variables of Q to constants such that by replacing every variable X by $G(X)$, every atom in Q becomes a fact in D . We denote this set of facts by $\text{facts}(Q, D, G)$. A grounding for a UCQ is a grounding for one of its CQs. Each grounding G yields an *output tuple*, which is the restriction of G to the head variables of Q . Multiple groundings may yield the same output tuple. We use $G(Q, D, \mathbf{t})$ to denote the set of groundings yielding \mathbf{t} . We use $Q(D)$ to denote the set of output tuples for Q w.r.t. D .

Query Lineage. Consider a database $D = D_n \cup D_x$. To construct the lineage of a query over D , we first associate each endogenous fact (or tuple) f in D_n with a Boolean variable $v(f)$. Given a UCQ Q and an output tuple \mathbf{t} , the lineage of \mathbf{t} with respect to Q over D , denoted by $\text{lin}(Q, D, \mathbf{t})$, is a positive Boolean formula in Disjunctive Normal Form over the variables $v(f)$ of facts f in D_n :

$$\text{lin}(Q, D, \mathbf{t}) = \bigvee_{G \in G(Q, D, \mathbf{t})} \bigwedge_{f \in \text{facts}(Q, D, G) \cap D_n} v(f)$$

That is, the lineage is a disjunction over all groundings yielding \mathbf{t} . For each such grounding G , the lineage has a conjunctive clause consisting of all endogenous facts in $\text{facts}(Q, D, G)$. If Q is a Boolean query, we abbreviate $\text{lin}(Q, D, ())$ by $\text{lin}(Q, D)$.

For an aggregate query $\langle \alpha, \gamma, Q \rangle$, let $Q(D) = \{t_1, \dots, t_k\}$. The lineage of the query over D is a semimodule expression of the form:

$$\text{lin}(\langle \alpha, \gamma, Q \rangle, D) = \sum_M \text{lin}(Q, D, t_i) \otimes \gamma(t_i),$$

where the summation uses the $+_M$ operator of the monoid M corresponding to α . For simplicity, we consider any fact f such that $v(f)$ is not in the lineage to be exogenous.

Example 2.8. Table 1 presents the lineage of the UCQ Q_1 and the aggregate query Q_2 in Figure 1. The lineage of Q_1 is a DNF boolean formula, while the lineage of Q_2 is a semimodule expression with the Max monoid. The boolean variables in the lineages correspond to tuples of the database in Figure 1. For instance, the variable a_1 corresponds to the Cast tuple 'Brad Pitt' annotated with a_1 .

Banzhaf and Shapley Values. We use the Banzhaf and Shapley values to measure the contribution of database facts to the answers of Boolean and aggregate queries [2, 19].

Definition 2.9 (Banzhaf Value). Given a Boolean or aggregate query Q and a database $D = D_x \cup D_n$, the Banzhaf value of an endogenous fact $f \in D_n$ is defined as:

$$\text{Banzhaf}(Q, f, D) \stackrel{\text{def}}{=} \sum_{D' \subseteq D_n \setminus \{f\}} Q(D' \cup \{f\} \cup D_x) - Q(D' \cup D_x)$$

Given a Boolean formula or a semimodule expression Ψ over a variable set X , the Banzhaf value of $x \in X$ is:

$$\text{Banzhaf}(\Psi, x) = \sum_{Y \subseteq X \setminus \{x\}} \Psi[Y \cup \{x\}] - \Psi[Y] \quad (1)$$

In case of a Boolean formula φ , the Banzhaf value of a variable x can be computed based on the model counts of the formulas obtained from φ by substituting x by the constants 1 and 0 [2]:

$$\text{Banzhaf}(\varphi, x) = \#\varphi[x := 1] - \#\varphi[x := 0] \quad (2)$$

Definition 2.10 (Shapley value). For sets Y, Z , let $C_Y^Z = C_Y = \frac{|Y|!(|X|-|Y|-1)!}{|X|!}$. Given a Boolean or aggregate query Q and a database $D = D_x \cup D_n$, the Shapley value of an endogenous fact $f \in D_n$ is defined as:

$$\text{Shapley}(Q, f, D) \stackrel{\text{def}}{=} \sum_{Y \subseteq D_n \setminus \{f\}} C_Y^{D_n} \cdot (Q(Y \cup \{f\} \cup D_x) - Q(Y \cup D_x))$$

Given a Boolean formula or a semimodule expression Ψ over a variable set X , the Shapley value of $x \in X$ is:

$$\text{Shapley}(\Psi, x) = \sum_{Y \subseteq X \setminus \{x\}} C_Y^X \cdot (\Psi[Y \cup \{x\}] - \Psi[Y]) \quad (3)$$

Prior work connected the Banzhaf and Shapley values for database facts and those of boolean variables [19]. Given a Boolean or aggregate query Q , a database $D = D_x \cup D_n$ and a fact $f \in D_n$:

$$\text{Banzhaf}(Q, D, f) = \text{Banzhaf}(\text{lin}(Q, D), v(f))$$

$$\text{Shapley}(Q, D, f) = \text{Shapley}(\text{lin}(Q, D), v(f))$$

Example 2.11. Consider the lineage of Q_1 , φ as presented in figure 1. Calculating the Banzhaf or Shapley value of the variable a_1 representing the *actor* fact of Brad Pitt. Consider the assignment $\theta = \{m_3, m_2\}$. The marginal contribution of a_1 on θ is $\varphi[\theta \cup \{a_1\}] - \varphi[\theta] = 1$ because θ is not a satisfying assignment for φ but $\theta \cup \{a_1\}$ is. The Banzhaf value is achieved by summing up the marginal contributions across all assignments. For Shapley value, the marginal contribution of a_1 to θ would be multiplied by a coefficient of $\frac{2!(7-2-1)!}{7!}$ since $|\theta| = 2$ and $|\text{vars}(\varphi)| = 7$. The Shapley value is achieved by summing up the marginal contributions across all assignments multiplied by the corresponding coefficients.

Decomposition Trees. Prior work uses knowledge compilation techniques to compute Banzhaf and Shapley values [2, 11, 17]. The key idea is to compile query lineage into a circuit that enables efficient computation of these values. We introduce decomposition trees for Boolean formulas and semimodule expressions. The definition we provide here differs slightly from prior work, as it incorporates non-binary gates [10].

Definition 2.12. A *decomposition tree*, or d-tree for short, is defined recursively as follows:

- If v is a 0/1 constant or a variable then v is a d-tree.

Algorithm 1 Lift

Input: Lifted DNF formula (φ, ℓ)

Output: Saturated lifted DNF formula

```

1: while  $(\varphi, \ell)$  is not saturated do
2:   if  $\varphi$  has a maximal set  $V$  of cofactor-equivalent variables
     with  $|V| > 1$  then
3:      $(\varphi, \ell) \leftarrow \text{lift-or}(\varphi, \ell, V)$ 
4:   if  $\varphi$  has a maximal set  $V$  of interchangeable variables with
      $|V| > 1$  then
5:      $(\varphi, \ell) \leftarrow \text{lift-and}(\varphi, \ell, V)$ 
6: return  $(\varphi, \ell)$ 
  
```

Example 3.2. Consider again the formulas φ_0 , φ_1 , and φ_2 from Example 3.1. The sets $\{x_1, x_2\}$ and $\{x_4, x_5\}$ are both maximal sets of interchangeable variables in φ_0 . The set $\{y_1, y_3\}$ is a maximal set of cofactor-equivalent variables in φ_1 with cofactor $\{y_2\}$. While the formulas φ_0 and φ_1 are not saturated, the formula φ_2 is.

Lifting DNF Formulas. We describe how a DNF formula can be systematically translated into a saturated lifted DNF formula. We first introduce some further notation. Consider a lifted DNF formula (φ, ℓ) over $X \cup Y$, where $\varphi = \bigwedge_{i \in [n]} C_i$. Let $V = \{y_1, \dots, y_k\}$ be a set of cofactor-equivalent variables in φ , where each variable has cofactor $\{N_1, \dots, N_p\}$. Let C_{i_1}, \dots, C_{i_m} be all clauses in φ containing variables from V . We define $\text{lift-or}(\varphi, \ell, V)$ to be the lifted DNF formula (φ', ℓ') , where: (1) φ' results from φ by omitting the clauses C_{i_1}, \dots, C_{i_m} and adding the new clause $y \wedge \bigwedge_{i \in [p]} N_i$ for a fresh variable $z \in Y \setminus \text{vars}(\varphi)$; (2) the function ℓ' is defined by $\ell'(y) = \bigvee_{i \in [k]} \ell(y_i)$ and $\ell'(y') = \ell(y')$ for all $\ell' \in \text{vars}(\varphi)$ with $y' \neq y$. Now, assume that $V = \{y_1, \dots, y_k\}$ consists of interchangeable variables in φ . We denote by $\text{lift-and}(\varphi, \ell, V)$ the lifted DNF formula (φ', ℓ') , where: (1) $\varphi' = C'_1 \wedge \dots \wedge C'_n$ such that, given a fresh variable $y \in Y \setminus \text{vars}(\varphi)$, each C'_i results from C_i by replacing the variables y_1, \dots, y_k by y ; (2) ℓ' is defined by $\ell'(y) = \bigwedge_{i \in [k]} \ell(y_i)$ and $\ell'(y') = \ell(y')$ for all $y' \in \text{vars}(\varphi)$ with $y' \neq y$.

The function **Lift** in Algorithm 1 transforms any lifted DNF formula φ into a saturated lifted DNF formula. It repeatedly calls the function **lift-or** or **lift-and** as long as φ contains a variable set V with $|V| > 1$, such that all variables in V are either cofactor-equivalent or interchangeable, respectively.

Proposition 3.3. Let φ be a DNF formula, ℓ the identity function on $\text{vars}(\varphi)$, and (φ', ℓ') the output of **Lift** (φ, ℓ) in Algorithm 1. Then, (φ', ℓ') is a saturated lifted formula such that: (1) $\text{inline}(\varphi', \ell')$ is equivalent to φ . (2) $\ell'(x)$ and $\ell'(y)$ are independent formulas for each distinct variables $x, y \in \text{vars}(\varphi')$. (3) $\ell'(x)$ is a read-once formula for each $x \in \text{vars}(\varphi')$.

Example 3.4. We explain how the procedure **Lift** transforms the lineage φ_0 of the query Q_1 depicted in Table 1 (top left). Let ℓ_0 be the identity function that maps every variable in φ_0 to itself. The variables d_1 and d_2 have the same cofactor $\{\{a_1, m_3\}, \{a_1, m_2\}, \{a_3, m_3\}, \{a_2, m_3\}, \{a_3, m_2\}\}$. The procedure executes **lift-or** $(\varphi_0, \ell_0, \{d_1, d_2\})$, which returns the lifted formula (φ_1, ℓ_1) , where

$$\varphi_1 = (d_{1,2} \wedge a_1 \wedge m_3) \vee (d_{1,2} \wedge a_1 \wedge m_2) \vee (d_{1,2} \wedge a_3 \wedge m_3) \vee (d_{1,2} \wedge a_2 \wedge m_3) \vee (d_{1,2} \wedge a_3 \wedge m_2),$$

and $\ell_1 = \{d_{1,2} \mapsto (d_1 \vee d_2)\} \cup \{x \mapsto x \mid x \in \text{vars}(\varphi) \text{ and } x \neq d_{1,2}\}$. In φ_1 , the variables a_1 and a_3 have the same cofactor $\{\{d_{1,2}, m_3\}, \{d_{1,2}, m_2\}\}$. The procedure executes **lift-or** $(\varphi_1, \ell_1, \{a_1, a_3\})$ and obtains (φ_2, ℓ_2) shown in Table 1 (top right). The formula φ_2 does not have two or more variables that are cofactor-equivalent or interchangeable. Hence, (φ_2, ℓ_2) is saturated.

Algorithm 2 LiftedCompile

Input: DNF lineage φ for query Q on database D

Output: d-tree for φ

```

1:  $(\psi, \ell) \leftarrow \text{Lift}(\varphi, \ell')$  where  $\ell'$  is identity function on  $\text{vars}(\varphi)$ 
2: return  $\text{Compile}(\psi, \ell)$ 
  
```

```

3: procedure  $\text{Compile}(\psi, \ell)$ 
4:   if  $\psi$  is a variable  $x$  then return  $\ell(x)$ 
   switch  $\psi$ 
5:     case  $\psi_1 \vee \dots \vee \psi_n$  for independent  $\psi_1, \dots, \psi_n$ 
6:        $T_\psi \leftarrow \bigoplus_{i \in [n]} \text{Compile}(\psi_i, \ell|_{\psi_i})$ 
7:     case  $\psi_1 \wedge \dots \wedge \psi_n$  for independent  $\psi_1, \dots, \psi_n$ 
8:        $T_\psi \leftarrow \bigodot_{i \in [n]} \text{Compile}(\psi_i, \ell|_{\psi_i})$ 
9:     default
10:      Pick a most frequent variable  $y$  in  $\psi$ 
11:       $T_\psi \leftarrow (y \odot \text{Compile}(\text{Lift}(\psi[y := 1], \ell|_{\psi[y:=1]}))) \oplus$ 
          $(\neg y \odot \text{Compile}(\text{Lift}(\psi[y := 0], \ell|_{\psi[y:=0]})))$ 
   return  $T_\psi$ 
  
```

D-Tree Compilation Using Lifting. Algorithm 2 gives our lifted compilation procedure for the query lineage (Boolean formula) into a d-tree. The input lineage φ is first translated into a saturated lifted formula (ψ, ℓ) (Line 1). The lifted formula is then passed to the sub-procedure **_Compile**, which traverses recursively over the structure of ψ . If ψ is a single variable x , the procedure returns $\ell(x)$, which is a read-once formula (Line 4). If ψ is a disjunction of independent sub-formulas ψ_1, \dots, ψ_n , it first constructs the d-trees for $(\psi_1, \ell|_{\psi_1}), \dots, (\psi_n, \ell|_{\psi_n})$, where $\ell|_{\psi_i}$ is the restriction of the domain of ℓ onto the variables of ψ_i . Then, it combines the d-trees using \bigoplus (Line 5). If ψ is a conjunction of independent sub-formulas ψ_1, \dots, ψ_n , the procedure combines the d-trees of the sub-formulas using \bigodot (Line 7). Otherwise, it performs Shannon expansion on a variable y as follows. It constructs d-trees T_0 for $\text{Lift}(\psi[y := 0], \ell|_{\psi[y:=0]})$ and T_1 for $\text{Lift}(\psi[y := 1], \ell|_{\psi[y:=1]})$. Then, it combines these d-trees using the $\sqcup_{\ell(y)}$ gate with 0-branch T_0 and 1-branch T_1 . (Lines 10-11). The lifting procedure is applied to $\psi[y := 1]$ and $\psi[y := 0]$, since the substitution of y by a constant can create new equivalences between other variables.

Example 3.5. We apply Algorithm 2 to the lineage of Q_1 in Table 1 (top left). The saturated lifted lineage (φ, ℓ) constructed in Line 1 is given in Table 1 (top right). Since the variable $d_{1,2}$ appears in all clauses of φ , we factor out $d_{1,2}$ from all clauses and decompose the formula using independent-and: $d_{1,2} \odot ((a_{1,3} \wedge m_3) \vee (a_{1,3} \wedge m_2) \vee (a_2 \wedge m_3))$ and proceeds recursively on both sides of \odot . Since ℓ maps $d_{1,2}$ to the read-once formula $d_1 \vee d_2$, we return the latter formula for $d_{1,2}$. The formula on the right side of \odot can only be decomposed using Shannon expansion on one of its variables. We choose $a_{1,3}$, one of its most frequent variables. We then perform lifting on each of the new formulas resulted from the Shannon expansion. This

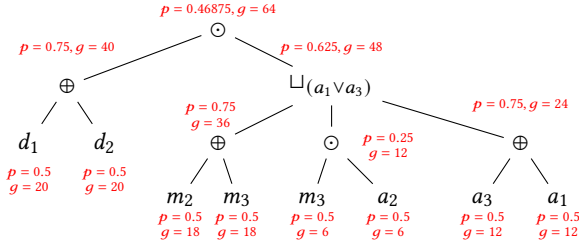


Figure 2: Lifted d-tree for the lineage in Figure 1. Each node is annotated with probability p and partial derivative g computed during the Banzhaf value computation with the gradient approach. Node m_3 is shown once for lack of space.

process continues until we obtain the d-tree from Figure 2 (ignore for now the numeric values appearing next to nodes).

Examples 3.4 and 3.5 demonstrate two benefits of our lifting approach over compilation without lifting: faster compilation time and smaller resulting d-tree. The compilation is more efficient since the lifted lineage can be small and Shannon expansion can be applied to the new fresh variables, which can represent large formulas.

3.2 Gradient-Based Computation

We propose a novel approach for Banzhaf and Shapley computation, yielding the following result:

Theorem 3.6. For a given d-tree T , the Banzhaf and Shapley values of all variables in T can be computed in time $O(|T|)$ and respectively $O(|T| \cdot |\text{vars}(T)|^2)$.

Our approach interprets a d-tree as a function over random variables and reduces the computation of Banzhaf and Shapley values to the evaluation and computation of the gradient of this function.² We first introduce our approach for the computation of Banzhaf values and then show how to extend it to Shapley values.

Probabilistic Interpretation of D-Trees. We associate with each Boolean variable x occurring in a d-tree leaf, a variable p_x taking values in $[0, 1]$. A value assigned to p_x is the probability for $x = 1$. $Pr[T]$ is then a function over these variables, defined recursively using the definitions in Table 2. Its value is the probability that the Boolean formula represented by the d-tree T evaluates to 1.

From $Pr[T]$ to Banzhaf values. We next connect the Banzhaf value of a variable x to the partial derivative of $Pr[T]$ w.r.t. x . The proof is deferred to the full version [3].

Proposition 3.7. Given a d-tree T and a variable $x \in \text{vars}(T)$, the following holds:

$$\text{Banzhaf}(T, x) = 2^{|\text{vars}(T)|-1} \cdot \left(\frac{\partial Pr[T]}{\partial p_x} \left(\frac{\vec{1}}{2} \right) \right)$$

where p_x is the Boolean variable corresponding to x , $\frac{\partial Pr[T]}{\partial p_x}$ is the partial derivative of $Pr[T]$ with respect to p_x , and $\frac{\vec{1}}{2}$ is the column vector $(\frac{1}{2}, \dots, \frac{1}{2})$ of length $|\text{vars}(T)|$.

²There are two seminal works related to our gradient-based approach: The Baur-Strassen result on the efficient computation of a multivariate function defined by an arithmetic circuit and its gradient vector [5], and the Darwiche result on efficient inference in Bayesian networks [9]. To the best of our knowledge, no prior work uses gradient-based computation for Banzhaf/Shapley values.

Table 2: Equations defining the probability and the partial derivatives for different gates. T is the d-tree rooted at the gate and the T_i 's are its child sub-trees.

Gate	Probability Expression $Pr[T]$	Partial Derivative
\oplus	$1 - \prod_{i \in [n]} (1 - Pr[T_i])$	$\frac{\partial(Pr[T])}{\partial(Pr[T_i])} = \prod_{j \in [n] \setminus \{i\}} (1 - Pr[T_j]) = \frac{1 - Pr[T]}{1 - Pr[T_i]}$
\odot	$\prod_{i \in [n]} Pr[T_i]$	$\frac{\partial(Pr[T])}{\partial(Pr[T_i])} = \prod_{j \in [n] \setminus \{i\}} Pr[T_j] = \frac{Pr[T]}{Pr[T_i]}$
\sqcup_f	$Pr[f] \cdot Pr[T_1] + (1 - Pr[f]) \cdot Pr[T_0]$	$\frac{\partial(Pr[T])}{\partial(Pr[T_1])} = Pr[f]$ $\frac{\partial(Pr[T])}{\partial(Pr[T_0])} = 1 - Pr[f]$ $\frac{\partial(Pr[T])}{\partial(Pr[f])} = Pr[T_1] - Pr[T_0]$

Efficient Computation. Table 2 gives the computation rules for the probability of a d-tree T based on the probabilities of its sub-trees; this follows standard probability theory. For each sub-tree T_i of T , we can also compute the partial derivative of $Pr[T]$ w.r.t. $Pr[T_i]$, where the latter is considered a variable in the former. Using the chain rule, we can then compute the gradient of $Pr[T]$ w.r.t. the variables at the leaves of T .

The procedure GradientBanzhaf in Algorithm 3 computes probabilities and partial derivatives at $\frac{\vec{1}}{2}$ (normalized by $2^{|\text{vars}(T)|-1}$) as follows. Given a d-tree T_ϕ , the algorithm gradually computes two quantities for each tree node v : $v.p$ which is the probability of T_v , and $v.g$ which is $\frac{\partial Pr(T_\phi)}{\partial Pr(T_v)}(\frac{\vec{1}}{2})$, where T_v is the d-tree rooted at v . First, $v.p$ values are initialized to $\frac{1}{2}$ and the $v.g$ value of the root is initialized to $2^{|\text{vars}(T_\phi)|-1}$ to account for normalization. Probabilities are computed in a bottom-up fashion, using the expressions in the second column of Table 2 (Lines 1 - 3). Partial derivatives are computed in a top-down fashion, using the expressions in the third column of Table 2 and the chain rule (Lines 5-6). For a variable v , the sum of partial derivatives for different leaf nodes representing v is equal to $\frac{\partial Pr[T]}{\partial p_v}(\frac{\vec{1}}{2})$ which equals its Banzhaf value by Prop. 3.7.

Example 3.8. Reconsider the d-tree from Figure 2, where its nodes are now annotated with p and g values. We explain the computation for the \oplus node v that is the left child of the root. Its p value is computed using the equation for \oplus from Table 2 and the probability of its sub-trees d_1 and d_2 : $v.p = 1 - (1 - d_1.p) \cdot (1 - d_2.p) = 1 - 0.5 \cdot 0.5 = 0.75$. Its g value is: $v.g = v.parent.g \cdot \frac{\partial v.parent.p}{\partial v.p} = v.parent.g \cdot \frac{v.parent.p}{v.p} = 64 \cdot \frac{0.46875}{0.75} = 40$ where we chose the equation for the derivative according to the gate of v 's parent.

From Banzhaf to Shapley values. Our approach can be adapted to Shapley values. The Shapley values are linked to probabilities and derivatives using the following equation:

$$\text{Shapley}(T, x) = 2^{|\text{vars}(T)|-1} \cdot \sum_{k \in [|\text{vars}(T)|]} C_k \cdot \left(\frac{\partial Pr_k[T]}{\partial p_x} \left(\frac{\vec{1}}{2} \right) \right)$$

where the k -probability $Pr_k[T] = \sum_{S \subseteq \text{vars}(T), |S|=k} Pr[S] \cdot T[S]$ is the probability of T computed only over the valuations of size k . To compute Shapley values, the equations in Table 2 are adapted to compute $Pr_k[T]$ instead of $Pr[T]$, and $\frac{\partial(\sum_k C_k \cdot Pr_k[T])}{\partial Pr_j[T_i]}(\frac{\vec{1}}{2})$ for each $j \in [n]$ instead of $\frac{\partial Pr[T]}{\partial Pr[T_i]}(\frac{\vec{1}}{2})$. Similar equations to those shown in Table 2 are obtained for these partial derivatives. Lines 1 - 3 are then changed to compute a vector of all k -probabilities for each node. Line 4 is modified to initialize the root's vector of derivatives to be

Algorithm 3 GradientBanzhaf

Input: D-tree T_φ for formula φ

Output: Banzhaf values for all variables in φ

- 1: Initialize $v.p \leftarrow \frac{1}{2}$ for each variable node v
 - 2: **for** each node v in the tree in bottom-up order **do**
 - 3: Compute $v.p$ according to Table 2 and $v.gate$
 - 4: $T.root.g = 2^{|\text{vars}(\varphi)|-1}$
 - 5: **for** each node v excluding $T.root$ in top-down order **do**
 - 6: $v.g \leftarrow v.parent.g \cdot \frac{\partial v.parent.p}{\partial v.p}$ according to Table 2 and $v.parent.gate$
 - 7: **for** each variable $x \in \text{vars}(\varphi)$ **do**
 - 8: $\text{Banzhaf}(T_\varphi, x) \leftarrow \sum_{\text{node } v \text{ in } T \text{ for variable } x} v.g$
 - 9: **return** Banzhaf(T_φ, x) for all variables $x \in \text{vars}(\varphi)$
-

the Shapley coefficients, and Lines 5 - 6 are modified to compute the partial derivative for each value of j .

Algorithms. The LEXABAN (LEXASHAP) algorithm for computing Banzhaf (Shapley) values uses Algorithm 2 to compile the query lineage to a d-tree T and then Algorithm 3 (its adaptation to Shapley values) to T to compute Banzhaf (Shapley) values for all variables in T , or equivalently for all tuples contributing to the lineage.

Complexity. For each node in the d-tree T , its probability can be computed in time linear in the number of its children, so in time $O(|T|)$ for all nodes in T . For Banzhaf computation, the gradient, i.e., the vector of partial derivatives w.r.t. each leaf node can be computed in $O(1)$ per node using the equations in Table 2, so in $O(|T|)$ time for the entire gradient. For Shapley computation, we need to compute the k -probabilities at each node v and their partial derivatives; the number of the k -probabilities is $|\text{vars}(T_v)|$. Each k -probability and partial derivative can be computed in $O(|\text{vars}(T)|)$ time, resulting in both parts of the algorithm performing $O(|\text{vars}(T)|^2)$ computations per node.

4 ATTRIBUTION FOR AGGREGATE QUERIES

In this section, we introduce algorithms for computing Banzhaf and Shapley values for aggregate queries based on d-trees for semimodule expressions (Section 2).

4.1 Basic Algorithms

We separately consider linear aggregates (SUM and COUNT) and non-linear ones (MIN and MAX).

SUM and COUNT. Let $\Phi = \sum_i \varphi_i \otimes m_i$ be a semimodule expression where the sum uses the $+_M$ -operator of the monoid, and let $\varphi = \bigvee \varphi_i$ be the Boolean part of Φ . In case the aggregate is SUM, $+_M$ is the arithmetic $+$, and we obtain:

$$\begin{aligned} \text{Banzhaf}(\Phi, x) &= \sum_{S \subseteq \text{vars}(\varphi) \setminus \{x\}} \Phi[S \cup \{x\}] - \Phi[S] \\ &= \sum_{S \subseteq \text{vars}(\varphi) \setminus \{x\}} \sum_i (\varphi_i[S \cup \{x\}] - \varphi_i[S]) \cdot m_i = \sum_i \text{Banzhaf}(\varphi_i, x) \cdot m_i \end{aligned}$$

Analogously, it holds that $\text{Shapley}(\Phi, x) = \sum_i \text{Shapley}(\varphi_i, x) \cdot m_i$ (note that this is analogous to the observation already made in [19]). For the COUNT aggregate, we set $m_i = 1$ for each m_i , and then the above equalities hold. Thus, to compute the Banzhaf/Shapley value

for a variable x in Φ , we (1) construct a d-tree T_i for each φ_i , (2) compute for each T_i and each $x \in \varphi_i$ the the Banzhaf/Shapley value for x as in Section 3, and (3) derive the Banzhaf/Shapley of x in Φ using the above equations.

MIN and MAX. For MIN and MAX which are not linear, we adopt a different approach. Namely, we first construct a d-tree for the full semimodule expression and then compute Banzhaf/Shapley values from the d-tree. The former was shown in [12] and we next show how the latter is performed.

Let $\text{PV}_\Phi = \{\Phi(\theta) \mid \theta \subseteq \text{vars}(\Phi)\}$ be the set of values to which the semimodule expression Φ can evaluate to under possible valuations. If Φ is the lineage of an aggregate query $\langle \text{MIN}, \gamma, Q \rangle$ or $\langle \text{MAX}, \gamma, Q \rangle$, the size of PV_Φ is bounded by the size of the result of Q , which is polynomial in the input database size. We next relate Banzhaf/Shapley values of a variable x in a semimodule expression with the model counts of the expression under substitutions that set x to 0 or 1:

Proposition 4.1. Given a semimodule expression Φ and a variable $x \in \text{vars}(\Phi)$, it holds:

$$\begin{aligned} \text{Banzhaf}(\Phi, x) &= \sum_{p \in \text{PV}_\Phi} (\#^p \Phi[x := 1] - \#^p \Phi[x := 0]) \cdot p \\ \text{Shapley}(\Phi, x) &= \sum_{p \in \text{PV}_\Phi, k \in [|\text{vars}(\Phi)|]} [(\#^p \Phi[x := 1] - \#^p \Phi[x := 0]) \cdot p] \cdot C_k \\ \text{where } C_k &= \frac{|k|! \cdot (|\text{vars}(\Phi)| - k - 1)!}{|\text{vars}(\Phi)|!}. \end{aligned}$$

We can further show, via a dynamic programming algorithm:

Proposition 4.2. Let T be a d-tree representing a semimodule expression $\Phi \in \text{PosBool}(X) \otimes M$ s.t. $M = (\mathbb{R}, \max, -\infty)$ or $M = (\mathbb{R}, \min, \infty)$. Let $p \in \text{PV}_\Phi$ and $k \in [|\text{vars}(T)|]$. We can compute $\#^p \Phi$ and $\#^p_k \Phi$ in PTIME.

Combined, the two propositions provide a PTIME algorithm for computing Banzhaf and Shapley values given a d-tree representation of the lineage of an aggregate query.

4.2 Lifting for Aggregates

We extend our technique of lifting variables presented in Section 3.1 from UCQs to aggregate queries. First, for linear aggregates (SUM and COUNT), our solution in Section 4.1 constructs a separate d-tree for each formula φ_i in the Boolean part $\bigvee_i \varphi_i$ of the given semimodule expression. Thus, we can apply the procedures Lift and LiftedCompile in Algorithms 1 and 2 to each formula φ_i . For MIN and MAX, new machinery is needed. In the following, we show the construction for MAX. The construction for MIN is analogous.

Extending the Definitions. Given disjoint variable sets X and Y , a *lifted semimodule expression* over $X \cup Y$ is a pair (φ, ℓ) where φ is a boolean formula over $X \cup Y$ and $\ell : \text{vars}(\varphi) \rightarrow \text{PosBool}(X) \cup (\text{PosBool}(X) \otimes M)$ is a mapping from variables in φ to either semimodule expressions or Boolean formulas. If φ is in DNF, we refer to (φ, ℓ) as a *lifted DNF semimodule expression*. The *inlining* $\text{inline}(\varphi, \ell)$ is obtained in two steps: First, we replace each variable $x \in \text{vars}(\varphi)$ by $\ell(x)$; then, we replace each \vee between semimodule expressions by $+_M$ and each \wedge between a formula and a semimodule expression by \otimes . We call a lifted semimodule expression *valid* if exactly one variable in each DNF caluse is mapped to a semimodule expression.

Translating Lineage to a Lifted DNF Semimodule Expression. Recall that the lineage of an aggregate query is captured by a semimodule expression $\Phi = \sum \varphi_i \otimes m_i$, where each $\varphi_i = \bigvee_j C_{i,j}$ is in DNF. We transform Φ to the equivalent form $\Phi' = \sum_{i,j} C_{i,j} \otimes m_i$. This preserves equivalence for $(\mathbb{R}, \max, -\infty)$ and $(\mathbb{R}, \min, \infty)$ due to their idempotence. We further transform Φ' into a lifted DNF semimodule expression (φ, ℓ) with $\varphi = \bigvee_{i,j} (C_{i,j} \wedge w_i)$ and $\ell = \{w_i \mapsto m_i\} \cup \{x \mapsto x \mid x \in \text{vars}(\varphi) \text{ and } x \neq w_i\}$. Note that (φ, ℓ) is valid and its inlining is equivalent to Φ . We add this translation as an initial preprocessing step to Algorithm 2.

Extending the Lifting Algorithm. Given a lifted DNF semimodule expression (φ, ℓ) , cofactor-equivalence and interchangeability are defined as in the case of lifted Boolean DNF formulas. The application of lift-and or lift-or to a set of cofactor-equivalent or respectively interchangeable variables maintains the validity of the expression, which can be verified using Definition 2.6. Algorithm 1 works without changes for lifted DNF semimodule expressions.

Example 4.3. Consider the semimodule expression that is the lineage of Query Q_2 in Table 1. We first translate the expression to the lifted DNF semimodule expression $\varphi = (a_1 \wedge m_3 \wedge t_1) \vee (a_2 \wedge m_3 \wedge t_1) \vee (a_3 \wedge m_3 \wedge t_1) \vee (a_1 \wedge m_2 \wedge t_2) \vee (a_3 \wedge m_2 \wedge t_2) \vee (a_4 \wedge m_1 \wedge t_3)$; $\ell = \{t_1 \mapsto 377, t_2 \mapsto 322, t_3 \mapsto 176\} \cup \{x \mapsto x \mid x \in \text{vars}(\varphi) \setminus \{t_1, t_2, t_3\}\}$. Then, we execute the initial lifting in Algorithm 2. The algorithm searches for cofactor-equivalent or interchangeable sets of variables. It detects that a_1 and a_3 are cofactor-equivalent, thus, it performs lift-or on them. It also detects that each of the sets $\{m_3, t_1\}$, $\{m_2, t_2\}$ and $\{a_4, m_1, t_3\}$ are interchangeable and thus performs lift-and on them. The resulting expression after these operations is given in Table 1. The algorithm continues compilation as described in Section 3.1 to achieve a d-tree for the lineage.

4.3 Gradients for Aggregates

We now adapt the gradient-based computation from Section 3.2 to the MAX aggregate (the case of MIN is similar). In contrast to d-trees with Boolean outcomes, a semimodule expression can yield a numeric value for a given valuation over its Boolean variables. Our probabilistic interpretation thus looks at *expected values*.

Expected Values and Gradients. Given d-tree T for s semimodule expression Φ , we denote by $E[T]$ the expected value of the d-tree given probabilities p_x for the boolean variables. That is: $E[T] = \sum_{S \subseteq \text{vars}(T)} \Pr[S] \cdot T[S]$ (where in contrast to Section 3.2, $T[S]$ is now not necessarily in $[0, 1]$). Proposition 3.7 then trivially extends to expected values:

$$\text{Banzhaf}(T, x) = 2^{|\text{vars}(T)|-1} \cdot \left(\frac{\partial E[T]}{\partial p_x} \left(\frac{1}{2} \right) \right) \quad (4)$$

We further extend the definition of probabilities and k-probabilities, so that $\Pr[T = m_i] = \sum_{S \subseteq \text{vars}(T), T[S]=m_i} \Pr[S]$ and $\Pr_k[T = m_i] = \sum_{S \subseteq \text{vars}(T), |S|=k, T[S]=m_i} \Pr[S]$.

Now, using the linearity of expectation, we can write:

$$\frac{\partial E[T]}{\partial p_x} = \sum_{m_i \in \text{PV}_\Phi} m_i \cdot \frac{\partial \Pr[T = m_i]}{\partial p_x}$$

From Gradients to Banzhaf and Shapley. Similarly to Section 3.2, we obtain Banzhaf and Shapley values of all facts via partial derivatives. Concretely, we introduce an order over the variables,

and then use variable names as indices (identifying each variable with its location in the order). Let J be the Jacobian matrix, namely $J_x^i = \frac{\partial (v_i \cdot \Pr[T=v_i])}{\partial p_x}$. Further let J^k be the matrix defined by $J_x^{k,i} = \frac{\partial (v_i \cdot \Pr_k[T=v_i])}{\partial p_x}$ for each variable x . We can show:

$$\begin{aligned} \text{Banzhaf}(T, x) &= 2^{|\text{vars}(T)|-1} \cdot \sum_{i \in [n]} J_x^i \left(\frac{1}{2} \right) \\ \text{Shapley}(T, x) &= 2^{|\text{vars}(T)|-1} \cdot \sum_{i \in [n]} \sum_{k \in [|\text{vars}(T)|]} J_x^{k,i} \left(\frac{1}{2} \right) \cdot C_k \end{aligned}$$

$$\text{where } C_k = \frac{|k|! \cdot |D_n - k - 1|!}{|D_n|!}.$$

Banzhaf and Shapley Computation. We adapt Algorithm 3 to compute Jacobians via back-propagation, for a d-tree representing the lineage of an aggregate query. The adaptations are as follows. Lines 1-3 are replaced with the computation of the probability for each different possible value. The root derivative in Line 4 is initialized to a vector of possible values. Lastly, the computation in Lines 5-6 is replaced with the computation of the Jacobian matrix values. This allows us to compute a vector/matrix of derivatives at each node and obtain the Banzhaf/Shapley values at the leaves according to Equation 4.

5 EXPERIMENTS

This section details our experimental setup and results.

5.1 Experimental Setup and Benchmarks

We implemented all algorithms in Python 3.11 and performed experiments on a Linux Debian 14.04 machine with 1TB of RAM and an Intel(R) Xeon(R) Gold 6252 CPU @ 2.10GHz processor.

Algorithms. We benchmarked variants of our algorithm LEXABAN, with and without the two optimization techniques, for SPJU and aggregate queries. We also compared to two state-of-the-art competitors that also use compilation to d-trees, but only work for SPJU queries and do not include the variable lifting and gradient optimizations [2]: EXABAN for exact computation and ADABAN for deterministic approximate computation with error guarantees. ADABAN was shown [2] to be successful for instances where exact algorithms failed and proved to be significantly faster and more accurate than a Monte-Carlo randomized approximation from prior work [11]. For Shapley value computation, we benchmarked our algorithm LEXASHAP and compared it to the state-of-the-art competitor EXASHAP [2] that is the Shapley variant of EXABAN. For aggregate queries with linear aggregates, we also compared to a variant where we use EXABAN as the algorithm of choice for each boolean formula and then use linearity as explained in Section 4.1 to compute Banzhaf values with respect to the aggregate queries.

Datasets. We tested the algorithms using 301 queries evaluated over three datasets: Academic, IMDB and TPC-H (SF1). The workload (see Table 3) is based on prior work on Banzhaf/Shapley values for query answering [2, 11]. The lineage for all output tuples of these queries was constructed using ProVSQLE [30]. We created two variants of each benchmark: one including SPJU queries, and one including aggregate queries. The SPJU variant includes all queries from IMDB and Academic as is (they do not include aggregate

Table 3: Statistics of the datasets used in the experiments.

Dataset	#Queries	# Lineages	#Vars avg/max	#Clauses avg/max
Academic	92	7,865	79 / 6,027	74 / 6,025
IMDB	197	986,030	25 / 27,993	15 / 13,800
TPC-H	12	165	1,918 / 139,095	863 / 75,983
TPC-H _{agg}	7	140	746 / 6959	233 / 2076

queries), and all queries without nested subqueries and with aggregates removed from TPC-H. The aggregate queries variant (named *Academic_{agg}*, *IMDB_{agg}*, *TPC – H_{agg}*) was constructed as follows. For TPC-H it includes 7 queries that were preserved as is, namely without removing aggregation. These are the queries for which ProVSQL managed to compute all lineages in an hour or less. To further extend the benchmark, we created multiple variants of each of these queries, by changing the aggregation function to MAX, COUNT and SUM (the solution for MIN is almost identical to that of MAX). For IMDB and Academic, we created a synthetic aggregate variant with each of MAX, COUNT and SUM introduced to each query, and applied over randomly chosen values. These datasets are highly challenging, including very large semimodule lineage expressions. Overall, the resulting dataset consists of nearly 1M Boolean lineage expressions and 429 semimodule lineage expressions for each aggregate.

Measurements. We define an instance as the computation of the Banzhaf or Shapley values for all variables in a lineage of an output tuple of a query over a dataset. We report the success rate of each algorithm over the instances, where failure is declared in case an algorithm did not terminate an instance within one hour. We also report statistics of runtimes across all instances (average, median, maximal runtime, and percentiles). The pX columns in the tables give the execution times for the X-th percentile of the instances.

5.2 Summary of Experimental Findings

We next overview our main findings.

(1) For queries without aggregates, LEXABAN consistently and significantly outperforms the state-of-the-art algorithm EXABAN for Banzhaf computation and even the ADABAN approximation algorithm. LEXABAN was consistently faster than these baselines, achieving a speedup of up to 3 orders of magnitude for the hardest instances. This speedup also leads to LEXABAN achieving success on more than 90% of the boolean lineage expressions where EXABAN failed, (and more than 99.98% of the boolean expressions overall). (2) LEXABAN was successful on aggregate queries: For SUM and COUNT, LEXABAN succeeded in less than 100th of a second on all TPC-H_{agg} dataset and achieved success on all instances of *ACADEMIC_{agg}* and over 95% of the instances of *IMDB_{agg}*, far surpassing the extension of EXABAN to these instances. The non-linear MAX aggregate is far more challenging, and no (variant of) prior work is applicable for MAX queries; yet LEXABAN was successful for more than 93% of the instances of TPC-H_{agg}. It was successful for over 78% and 61% of the MAX queries from the highly challenging *ACADEMIC_{agg}* and *IMDB_{agg}* synthetic datasets respectively.

(3) The lifting technique significantly improved the compilation used by all variants of our approach: compilation time improved by more than 2 orders of magnitude compared to compiling without

Table 4: Query success rate: Percentage of queries for which the algorithms finish for all instances of a query. Lineage success rate: Percentage of instances (over all queries) for which the algorithms finish. The timeout is one hour.

Dataset	Algorithm	Query Success Rate	Lineage Success Rate
Academic	LEXABAN	100.00%	100.00%
	EXABAN	98.85%	99.99%
	ADABAN	98.85%	99.99%
IMDB	LEXABAN	95.43%	99.98%
	EXABAN	81.73%	99.63%
	ADABAN	87.82%	99.81%
TPC-H	LEXABAN	75.00%	93.33%
	EXABAN	58.33%	91.52%
	ADABAN	75.00%	92.73%

the lifting technique. In addition, compilation resulted in 1-2 orders of magnitude smaller d-trees.

(4) The gradient technique proved to be highly effective for large instances, pushing the boundary of computable instances: The technique yielded speedups across all instance sizes, with improvements exceeding two orders of magnitude for the largest instances.

(5) The lifting and gradient techniques also improve Shapley value computation: LEXASHAP shows speed-ups of over two orders of magnitude over EXASHAP. LEXASHAP achieves success on over 75% of the Boolean expression instances EXASHAP failed on.

5.3 Banzhaf Computation for SPJU Queries

Success rate. Table 4 shows the success rates of LEXABAN, EXABAN, and ADABAN, in the setting for which EXABAN and ADABAN were designed, namely queries without aggregates. LEXABAN achieves higher success rates at both lineage and query levels for all datasets. On Academic, LEXABAN succeeded on all instances. On IMDB, LEXABAN achieves 99.98% success at the lineage level, and more than 95% success at query level. LEXABAN succeeds on more than 95% of instances on which EXABAN failed, and 90% of instances on which ADABAN failed. An improvement (albeit of a lesser extent) was also observed for TPC-H.

Runtime performance. Table 5 gives the runtime of LEXABAN, EXABAN and ADABAN for the Academic, IMDB and TPC-H datasets. LEXABAN outperforms both EXABAN and ADABAN significantly. Despite achieving success for more instances, it keeps a smaller mean runtime compared with EXABAN and ADABAN. For Academic, except for the unique instance where it achieves success and the other algorithms fail on, runtime is below 5 seconds, which represents 119X improvement over EXABAN, and 36X improvement over ADABAN. Similar trends are observed for the IMDB and TPC-H datasets, where LEXABAN achieves 110X and respectively 106X improvements for the most expensive instance for EXABAN (Table 6). For the instances on which EXABAN failed (Table 7), LEXABAN still achieves success and low runtimes. For IMDB, this is an almost 95% success rate with a median runtime of 4.73 seconds, representing at least 761X speedup for more than half of those instances. For TPC-H, LEXABAN succeeded on 3 more instances; for one such instance LEXABAN needed 1.95 seconds, which means a speedup of at least 1846X over EXABAN which did not terminate in 1 hour.

Effect of lineage structure on runtime performance. Fig. 3 shows the effect of the number of variables and clauses in the query lineage on the success rates and runtimes of LEXABAN and EXABAN. When

Table 5: Runtime of LEXABAN, EXABAN and ADABAN. Note that the instance for Academic on which LEXABAN is slowest is one on which EXABAN and ADABAN failed on.

Dataset	Algorithm	Success rate	Execution times [sec]					
			Mean	p50	p90	p95	p99	Max
Academic	LEXABAN	100.00%	0.44	4E-4	2E-3	5E-3	0.14	3455
	EXABAN	99.99%	2.07	1E-3	0.01	0.20	164.5	563.5
	ADABAN	99.99%	0.76	1E-3	7E-3	0.05	60.05	173.7
IMDB	LEXABAN	99.98%	0.18	1E-3	5E-3	0.01	0.09	1219
	EXABAN	99.63%	1.58	2E-3	0.02	0.08	10.37	1793
	ADABAN	99.81%	1.82	1E-3	0.01	0.05	7.81	1802
TPC-H	LEXABAN	93.33%	0.12	7E-4	3E-3	0.62	1.95	8.73
	EXABAN	91.52%	4.23	0.89	0.94	51.05	61.98	69.18
	ADABAN	92.73%	2.37	3E-3	0.14	3.15	81.55	166.3

Table 6: Runtime of LEXABAN and EXABAN on instances for which EXABAN succeeded.

Dataset	algorithm	Execution times [sec]						
		Mean	p50	p75	p90	p95	p99	Max
Academic	LEXABAN	4E-3	4E-4	1E-3	2E-3	5E-3	0.14	4.99
	EXABAN	2.07	1E-3	2E-3	0.01	0.20	164.5	563.5
IMDB	LEXABAN	7E-3	1E-3	2E-3	4E-3	9E-3	0.05	16.26
	EXABAN	1.58	2E-3	3E-3	0.02	0.08	10.37	1793
TPC-H	LEXABAN	0.04	7E-4	8E-4	1E-3	0.53	0.65	0.65
	EXABAN	4.23	0.89	0.93	0.94	51.05	61.98	69.18

Table 7: Runtime of LEXABAN on instances for which EXABAN failed. There was one such instance in Academic.

Dataset	Success rate	Mean	p50	p75	p90	p95	p99	Max
Academic	100%	3455	-	-	-	-	-	3455
IMDB	94.74%	48.97	4.73	22.88	168.2	375.9	448.2	1219
TPC-H	21.43%	4.21	1.95	5.34	7.38	8.05	8.60	8.73

Table 8: Success rate and runtime of LEXABAN for the TPC – H_{agg} aggregate queries. Max(Naive) refers to running LEXABAN without optimizations on queries with MAX.

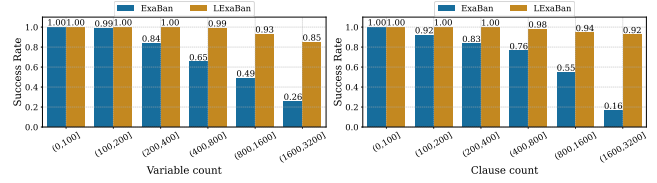
Aggregate	Success rate	Execution times [sec]						
		Mean	p50	p75	p90	p95	p99	Max
SUM	100%	2E-4	2E-6	2E-6	1E-3	2E-3	2E-3	2E-3
COUNT	100%	2E-4	1E-6	2E-6	1E-3	2E-3	2E-3	2E-3
MAX	93.57%	224.8	2E-3	3E-3	5E-3	2990	3349	3517
Max(Naive)	86.43%	3E-3	1E-3	1E-3	2E-3	3E-3	9E-3	0.26

the number of variables/clauses increases, LEXABAN’s success rate remains higher, while its runtime increases slower than for EXABAN. This highlights the advantages of both our techniques: lifting helps LEXABAN compile the lineage faster and into a smaller d-tree, while the gradient helps LEXABAN share the computation among many variables (see further analysis below).

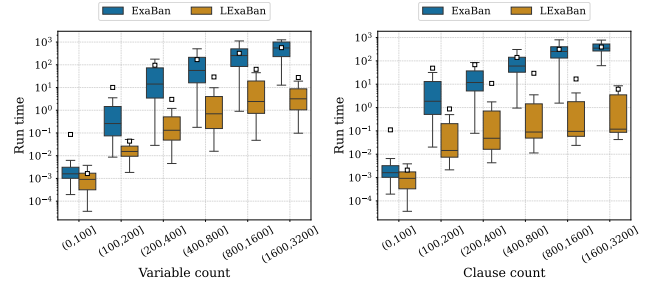
5.4 Banzhaf Computation for Aggregate Queries

We then tested LEXABAN on aggregate queries. Table 8 shows the results for the TPC-H_{agg} dataset. The performance is near-perfect for SUM and COUNT aggregates (incurring less than a millisecond for all instances). LEXABAN also terminated quickly for most instances with MAX aggregate; specifically, it significantly outperforms a naive variant of LEXABAN that does not include the optimizations.

For the synthetic and highly challenging IMDB_{agg} and Academic_{agg} datasets, we again observe good performance for SUM (and COUNT, for which the results were very similar and are omitted for lack of



(a) Success rate (average over all instances in each group)



(b) Runtime (ranges over all instances in each group)

Figure 3: Success rate and runtime of LEXABAN and EXABAN for all queries, grouped by the number of variables or clauses in the lineage. An interval $[i, j]$ represents the set of lineages with #vars (# clauses) between i and j .

Table 9: Success rate and runtime of LEXABAN and EXABAN for queries with SUM aggregate (similar results for COUNT aggregate are not shown) and MAX aggregates.

Dataset	Algorithm & Aggregate	Success rate	Execution times [sec]					
			Mean	p50	p90	p95	p99	Max
Academic _{agg}	LEXABAN _{sum}	100.00%	40.1	0.08	0.48	0.58	496	3561
	EXABAN _{sum}	97.70%	6.20	0.09	1.88	12.6	105	379.4
	LEXABAN _{max}	78.26%	15.11	3E-3	37.04	104.5	253.4	289.1
IMDB _{agg}	LEXABAN _{sum}	95.43%	37.54	0.06	9.77	71.56	774	2450
	EXABAN _{sum}	81.19%	62.7	0.27	65.8	117	1767	3286
	LEXABAN _{max}	61.93%	52.25	0.09	23.79	242.6	734.9	2407

space). Specifically, LEXABAN has succeeded on all SUM query instances for Academic_{agg} and 95.43% of such instances for IMDB_{agg}. When compared to a variant that uses EXABAN for the individual boolean formulas occurring in the semimodule expressions, LEXABAN was clearly superior. For MAX, there is no such baseline, and LEXABAN was still successful for over 78% for Academic_{agg} and over 61% for IMDB_{agg}. Furthermore, for both datasets, it terminated in less than 0.1 of a second for over 50 % of the lineages.

To better understand the performance over hard instances, we designed an experiment where we control the size of semimodule expressions by considering subsets of the tuples in the TPC-H Lineitem relation. Fig. 4 shows the average runtime for all lineages of a representative “hard” query (TPC-H Query 5). LEXABAN is able to compute the Banzhaf value on average in about 1.6 seconds, when the number of Lineitem facts is 3 million, while without the lifting and gradient optimizations (dubbed Naive in the figure), computation incurred on average over 10 minutes.

5.5 Effect of Lifting and Gradient Techniques

Fig. 5 shows that the variable lifting technique can speed-up the compilation of query lineage by over one order of magnitude on

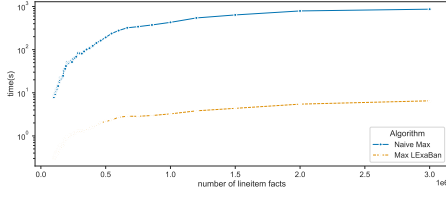
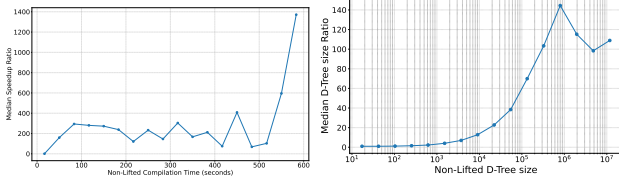


Figure 4: Runtime of LEXABAN for a TPC-H query 5 with MAX aggregate, when varying the number of Lineitem facts.



(a) D-tree compilation time with and without variable lifting. (b) D-tree size with and without variable lifting.

Figure 5: Effect of variable lifting on LEXABAN's compilation of lineage of Boolean queries into d-trees.

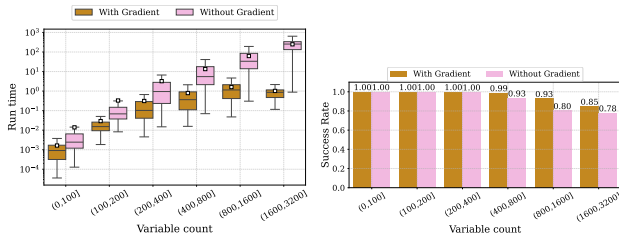


Figure 6: Effect of the gradient technique on the success rate and runtime of LEXABAN for lineage of Boolean queries.

average and reduce the size of the constructed d-trees by more than two orders of magnitude. These benefits increase as we consider increasingly larger lineage expressions. This experiment considers the lineage of all Boolean (Academic, IMDB, TPC-H) queries that can be processed within 10 minutes.

Fig. 6 shows the success rate and runtime of LEXABAN with and without the gradient technique for the lineage of all Boolean (Academic, IMDB, TPC-H) queries. The technique was beneficial for all considered instances, and its benefit increased with the instance size. For instances with more than 1600 variables, the average speedup is more than 2 orders of magnitude. One might expect a speedup of at most n for instances with n variables, since the gradient technique allows to compute the Banzhaf value simultaneously for all n variables. However, not all branches of the d-tree participate in the Banzhaf computation for each variable, and thus there is potential for significant computation sharing across variables.

5.6 Shapley Computation

We also report on the success rate and runtime of LEXASHAP, our extension of LEXABAN to compute the Shapley value. The main

Table 10: Runtimes for LEXASHAP and EXASHAP on instances for which EXASHAP succeeds.

Dataset	Algorithm	Execution times [sec]						
		Mean	p50	p75	p90	p95	p99	Max
Academic	LEXASHAP	0.02	2E-3	4E-3	9E-3	0.01	0.05	48.40
	EXASHAP	0.45	1E-3	3E-3	0.04	0.22	0.37	1397
IMDB	LEXASHAP	0.07	6E-3	0.02	0.05	0.10	0.57	520.3
	EXASHAP	3.35	3E-3	0.01	0.12	0.79	44.00	3574
TPCH	LEXASHAP	2E-3	2E-3	4E-3	5E-3	5E-3	5E-3	5E-3
	EXASHAP	7E-3	3E-3	3E-3	5E-3	9E-3	0.02	0.56

Table 11: Runtimes for LEXASHAP on instances for which EXASHAP fails.

Dataset	Success rate	Execution times [sec]						
		Mean	p50	p75	p90	p95	p99	Max
IMDB	79.13%	61.54	13.04	59.41	171.8	295.7	689.3	2835
TPCH	41.68%	22.90	23.94	26.82	27.89	29.02	29.92	30.14

competitor is EXASHAP from prior work [2], which also uses d-tree compilation but does not use our lifting and gradient techniques.

Table 10 shows the runtimes of LEXASHAP and EXASHAP for instances on which the latter succeeds. Table 11 shows the success rates and runtimes of LEXASHAP for instances on which EXASHAP fails. On average, LEXASHAP's speedup over EXASHAP is more than one order of magnitude across all datasets. LEXASHAP also has a much higher success rate: 79% for the IMDB instances on which EXASHAP fails, with a median time of 13.04 seconds; and 41% for the TPC-H instances on which EXASHAP fails, with a median runtime of 23.94 seconds. This represents a speedup of at least 276X and respectively 150X.

6 CONCLUSION

We have presented in this paper a novel practical solution for computing Banzhaf and Shapley values in query answering, thereby explaining query results. The solution is both more general and more efficient than the state-of-the-art: its generality is manifested in its support of aggregate queries; its efficiency is due to the incorporation of two novel techniques, namely *variable lifting* in the d-tree representation of query results lineage, and a *gradient-based approach* to efficiently compute Banzhaf and Shapley values based on the d-tree. Our experimental evaluation shows that the solution is faster than the state-of-the-art by several orders of magnitude. Future work includes extending our approach to broader query classes (e.g. negation) and developing efficient approximation schemes for Banzhaf and Shapley values, using our insights and techniques. Another line of work is exploring the applicability of our algorithms for other attribution measures such as the SHAP score.

REFERENCES

- [1] Serge Abiteboul, Richard Hull, and Victor Vianu. 1995. *Foundations of Databases*. Vol. 8. Addison-Wesley Reading. <http://webdam.inria.fr/Alice/>
- [2] Omer Abramovich, Daniel Deutch, Nave Frost, Ahmet Kara, and Dan Olteanu. 2024. Banzhaf Values for Facts in Query Answering. *Proc. ACM Manag. Data* 2, 3, Article 123 (2024), 26 pages. <https://doi.org/10.1145/3654926>
- [3] Omer Abramovich, Daniel Deutch, Nave Frost, Ahmet Kara, and Dan Olteanu. 2025. LEXABan-LEXASHap. <https://github.com/Omer-Abramovich/LEXABan-LEXASHap>. Accessed: 2025-03-01.
- [4] J.F. Banzhaf. 1965. Weighted voting doesn't work: A mathematical analysis. *Rutgers Law Review* 19, 2 (1965), 317–343.

- [5] Walter Baur and Volker Strassen. 1983. The complexity of partial derivatives. *Theoretical Computer Science* 22, 3 (1983), 317–330. [https://doi.org/10.1016/0304-3975\(83\)90110-X](https://doi.org/10.1016/0304-3975(83)90110-X)
- [6] Leopoldo Bertossi, Benny Kimelfeld, Ester Livshits, and Mikaël Monet. 2023. The Shapley Value in Database Management. *SIGMOD Rec.* 52, 2 (Aug. 2023), 6–17. <https://doi.org/10.1145/3615952.3615954>
- [7] Meghyn Bienvenu, Diego Figueira, and Pierre Lafourcade. 2024. Shapley value computation in ontology-mediated query answering. In *Proceedings of the 21st International Conference on Principles of Knowledge Representation and Reasoning (Hanoi, Vietnam) (KR '24)*. Article 15, 11 pages. <https://doi.org/10.24963/kr.2024/15>
- [8] James Cheney, Laura Chiticariu, and Wang-Chiew Tan. 2009. Provenance in Databases: Why, How, and Where. *Found. Trends Databases* 1, 4 (April 2009), 379–474. <https://doi.org/10.1561/19000000006>
- [9] Adnan Darwiche. 2003. A differential approach to inference in Bayesian networks. *J. ACM* 50, 3 (2003), 280–305. <https://doi.org/10.1145/765568.765570>
- [10] Adnan Darwiche and Pierre Marquis. 2002. A knowledge compilation map. *J. Artif. Int. Res.* 17, 1 (Sept. 2002), 229–264.
- [11] Daniel Deutch, Nave Frost, Benny Kimelfeld, and Mikaël Monet. 2022. Computing the Shapley Value of Facts in Query Answering. In *SIGMOD*. 1570–1583. <https://doi.org/10.1145/3514221.3517912>
- [12] Robert Fink, Larisa Han, and Dan Olteanu. 2012. Aggregation in probabilistic databases via knowledge compilation. *Proc. VLDB Endow.* 5, 5 (Jan. 2012), 490–501. <https://doi.org/10.14778/2140436.2140445>
- [13] Daniel Fryer, Inga Strümke, and Hien Nguyen. 2021. Shapley values for feature selection: The good, the bad, and the axioms. *Ieee Access* 9 (2021), 144352–144360.
- [14] Boris Glavic, Alexandra Meliou, and Sudeepa Roy. 2021. Trends in Explanations: Understanding and Debugging Data-driven Systems. *Found. Trends Databases* 11, 3 (Aug. 2021), 226–318. <https://doi.org/10.1561/19000000074>
- [15] Jonathan S. Golan. 1999. Semirings and their applications. <https://api.semanticscholar.org/CorpusID:122727231>
- [16] Ahmet Kara, Dan Olteanu, and Dan Suciu. 2024. From Shapley Value to Model Counting and Back. *Proc. ACM Manag. Data* 2, 2, Article 79 (May 2024), 23 pages. <https://doi.org/10.1145/3651142>
- [17] Pratik Karmakar, Mikaël Monet, Pierre Senellart, and Stephane Bressan. 2024. Expected Shapley-Like Scores of Boolean functions: Complexity and Applications to Probabilistic Databases. *Proc. ACM Manag. Data* 2, 2, Article 92 (2024), 26 pages. <https://doi.org/10.1145/3651593>
- [18] Majd Khalil and Benny Kimelfeld. 2022. The complexity of the Shapley value for regular path queries. *arXiv preprint arXiv:2212.07720* (2022).
- [19] Ester Livshits, Leopoldo Bertossi, Benny Kimelfeld, and Moshe Sebag. 2021. The Shapley Value of Tuples in Query Answering. *Logical Methods in Computer Science* Volume 17, Issue 3, Article 22 (Sep 2021). [https://doi.org/10.46298/lmcs-17\(3:22\)2021](https://doi.org/10.46298/lmcs-17(3:22)2021)
- [20] Scott M Lundberg and Su-In Lee. 2017. A unified approach to interpreting model predictions. In *NeurIPS*. 4765–4774. <http://papers.nips.cc/paper/7062-a-unified-approach-to-interpreting-model-predictions.pdf>
- [21] Xuan Luo, Jian Pei, Cheng Xu, Wenjie Zhang, and Jianliang Xu. 2024. Fast Shapley Value Computation in Data Assemblage Tasks as Cooperative Simple Games. *Proc. ACM Manag. Data* 2, 1, Article 56 (March 2024), 28 pages. <https://doi.org/10.1145/3639311>
- [22] Neha Makhija and Wolfgang Gatterbauer. 2024. Minimally Factorizing the Provenance of Self-join Free Conjunctive Queries. *Proc. ACM Manag. Data* 2, 2 (2024), 104. <https://doi.org/10.1145/3651605>
- [23] Alexandra Meliou, Wolfgang Gatterbauer, Katherine F. Moore, and Dan Suciu. 2010. The complexity of causality and responsibility for query answers and non-answers. *Proc. of VLDB Endow. (PVLDB)* 4, 1 (2010), 34–45. <https://www.vldb.org/pvldb/vol4/p34-meliou.pdf>
- [24] Alexandra Meliou, Wolfgang Gatterbauer, and Dan Suciu. 2011. Bringing Provenance to Its Full Potential Using Causal Reasoning. In *TaPP*, Peter Buneman and Juliana Freire (Eds.). USENIX Association.
- [25] Alexandra Meliou, Sudeepa Roy, and Dan Suciu. 2014. Causality and explanations in databases. *Proc. VLDB Endow.* 7, 13 (Aug. 2014), 1715–1716. <https://doi.org/10.14778/2733004.2733070>
- [26] Dan Olteanu and Jakub Zavodny. 2012. Factorised representations of query results: size bounds and readability. In *15th International Conference on Database Theory, ICDT '12, Berlin, Germany, March 26-29, 2012*, Alin Deutsch (Ed.). ACM, 285–298. <https://doi.org/10.1145/2274576.2274607>
- [27] L. S. Penrose. 1946. The Elementary Statistics of Majority Voting. *Journal of the Royal Statistical Society* 109, 1 (1946), 53–57. <http://www.jstor.org/stable/2981392>
- [28] Alon Reshef, Benny Kimelfeld, and Ester Livshits. 2020. The impact of negation on the complexity of the Shapley value in conjunctive queries. In *PODS*. 285–297.
- [29] Benedek Rozemberczki, Lauren Watson, Péter Bayer, Hao-Tsung Yang, Olivér Kiss, Sebastian Nilsson, and Rik Sarkar. 2022. The shapley value in machine learning. *arXiv preprint arXiv:2202.05594* (2022).
- [30] Pierre Senellart, Louis Jachiet, Silviu Maniu, and Yann Ramusat. 2018. Provsq: Provenance and probability management in postgresql. *Proc. of VLDB Endow. (PVLDB)* 11, 12 (2018), 2034–2037. <https://hal.inria.fr/hal-01851538/file/p976-senellart.pdf>
- [31] Lloyd S Shapley. 1953. A value for n-person games. *Contributions to the Theory of Games* 2, 28 (1953), 307–317. <http://www.library.fu.ru/files/Roth2.pdf#page=39>
- [32] Jianyuan Sun, Guoqiang Zhong, Kaizhu Huang, and Junyu Dong. 2018. Banzhaf random forests: Cooperative game theory based random forests with consistency. *Neural Networks* 106 (2018), 20–29.
- [33] Jiachen T Wang and Ruoxi Jia. 2023. Data banzhaf: A robust data valuation framework for machine learning. In *International Conference on Artificial Intelligence and Statistics*. PMLR, 6388–6421.