

Shredding higher-order nested queries

James Cheney

The University of Edinburgh
jcheney@inf.ed.ac.uk

Sam Lindley

The University of Edinburgh
Sam.Lindley@ed.ac.uk

Philip Wadler

The University of Edinburgh
wadler@inf.ed.ac.uk

Abstract

Reconciling high-level functional programming abstractions with the capabilities of databases is a major challenge. We present a modular account of *shredding*, the simulation of a single nested relational query by a number of flat relational queries. Our key insight is that shredding can be greatly simplified by first rewriting the input query into a canonical normal form. Normalisation allows us to define shredding translations on types and terms independently of one another. An added benefit of normalisation is that we get higher-order terms for free, provided that the result type is a plain nested relation type (without higher order components). In order to generate SQL we consider several alternatives for generating indexes, focusing on a lightweight use of SQL OLAP features.

We prove correctness of our translations, focusing on the central shredding step: shredding a nested query, running the shredded queries, and stitching the results back together yields the same results as running the nested query directly.

1. Introduction

Databases are arguably one of the most important applications of functional or declarative programming techniques. However, relational databases only support queries to flat tables, while programming languages typically provide complex data structures that allow arbitrary combinations of types including nesting one type constructor inside another. Motivated by this so-called *impedance mismatch*, and inspired by insights into language design based on monadic comprehensions [27], database researchers introduced nested relational query languages [5, 21] as a generalisation of flat relational queries to allow nesting collection types inside records or other types. Several recent language designs, such as XQuery [30] and PigLatin [17], have further extended these ideas, and they have been particularly influential on language-integrated querying systems such as Kleisli [29], Microsoft’s LINQ [15], Links [8], Ferry [10], Ur/Web [6], and SML# [16]. Of these, Links, Ferry and LINQ support queries over nested data structures (e.g. records containing sets) in principle, but in practice either reject such queries at run time or execute them inefficiently by fetching a large amount of data from the database and then executing the query in-memory.

This paper considers the problem of translating nested queries over nested data to flat queries over a flat representation of nested data, or *query shredding* for short. Our motivation is to support a free combination of the features of nested relational query lan-

guages with those of high-level functional programming languages. Query shredding has long been known to be possible in principle for pure nested relational query languages, as shown by Van den Bussche [25], and recently Grust et al. [10, 11] have developed a higher-order, nested relational query language called Ferry whose implementation performs shredding via a technique called *loop-lifting*. However, loop-lifting is a monolithic program transformation that is difficult to verify, and produces queries that make heavy use of advanced On-Line Analytic Processing (OLAP) features of SQL:1999 such as `row_number` that are challenging to optimise.

In this paper, we introduce a new, modular approach to query shredding. Our work is oriented towards the programming model of the Links programming language, but should be applicable to other systems based on similar ideas, including Ferry. Essentially, our idea is to decompose the translation from nested to flat queries into a number of simpler translations. Small, modular compilation passes are much easier to implement, verify, debug, and extend than a monolithic whole-program compilation pass. Our approach, like Ferry, employs some of SQL:1999’s OLAP features, but we delay their use as long as possible, so that at most one such operation is needed per generated query, and we identify cases where their use can be avoided.

Our approach works as follows. The first phase, *normalisation*, translates a higher-order query to a *normal form* in which higher-order features have been eliminated. The second phase, *shredding*, translates a single, nested query to a number of flat queries. These queries are organised in a *shredded package*, which is essentially a type expression whose collection type constructors are annotated with queries. The different queries are linked by *indexes*. The third phase, *let-insertion*, hoists correlated nested subqueries using a *let-binding* construct (equivalent to SQL’s `with`) and a row-numbering operation (equivalent to a restriction of SQL’s `row_number`). The fourth phase, *record flattening*, flattens nested records, and translation to SQL syntax is then trivial.

The normalisation phase builds on the rewriting approach taken in a series of papers by Wong [28], Cooper [7], and our prior work [13], but differs in some minor details (given in Appendix A) that simplify the proof of correctness and facilitate later stages. The shredding phase, and its proof of correctness, is new, and is the main contribution of the paper; it leverages the normalisation phase in that we can define translations on types and terms independently (in contrast to van den Bussche’s or Ferry’s approaches). The *let-insertion* phase is conceptually straightforward, but provides a vital link to proper SQL by providing an implementation of abstract indexes. Record flattening is standard, and the details are included in an appendix for completeness.

In contrast to van den Bussche’s simulation, our approach handles higher-order features and compiles to SQL. In contrast to Ferry, we do not just describe an algorithm but provide a proof of correctness. Our approach does have some limitations, which are the subject of current work: at present, it does not recognise idiomatic SQL-style grouping and aggregation, nor does it handle

queries that return functions. Some work has been done to handle both features in Ferry, albeit without any correctness proof [24].

Also, it is important to note that Ferry’s approach is geared towards a list-based interpretation of database queries, while we assume a bag-based semantics (matching proper SQL). This difference is easily overcome, as discussed in Section 8. In fact, the core technical part of the paper (Sections 4–6) works just as well for a list semantics. The only parts that rely on bag semantics are normalisation (Section 3) and conversion to SQL (Section 7).

1.1 An example

To motivate and illustrate our work, we present a small example showing how our shredding algorithm could be used to provide useful functionality to programmers working in a database programming language such as Ferry or Links. We first describe the code the programmer would actually write and the results the system produces. We then show how the shredding algorithm in this paper (combined with other, simpler transformations or techniques in previous work [7, 13]) implements this behaviour efficiently when the nested data is stored in tables on a relational database.

What the programmer sees Nested database schemas (i.e., type declarations) allow free mixing of collection (bag) types with record or base types. Consider a nested database schema for an organisation consisting of a collection of departments. Each department has a name, a collection of employees, and a collection of external contacts. Each employee has a name, salary and a collection of tasks. Some contacts are clients.

```

Task = String
Employee = ⟨name : String, tasks : Bag Task, salary : Int⟩
Contact = ⟨name : String, client : Bool⟩
Department = ⟨name : String, employees : Bag Employee,
              contacts : Bag Contact⟩
Organisation = Bag Department

```

The nested schema Σ contains a single collection *organisation* with type *Organisation*.

To illustrate the main features of shredding, we construct a query with two levels of nesting and a union operation. Suppose we wish to find for each department a collection of people of interest, both employees and contacts, along with a list of the tasks they perform. Specifically, we are interested in those employees that earn less than a thousand euros and those who earn more than a million euros, call them *outliers*, along with those contacts who are clients. The following code defines *poor*, *rich*, *outliers*, and *clients*:

```

isPoor x = x.salary < 1000
isRich x = x.salary > 1000000
filter p xs = for (x ← xs) where (p(x)) return x
outliers xs = filter (λx. isRich x ∨ isPoor x) xs
clients xs = filter (λx. x.client) xs
get xs f = for (x ← xs) return ⟨name = x.name, tasks = f x⟩

```

In our query language comprehensions are primitive, and we can define and use a higher-order *filter* function to build queries. Likewise, we introduce a higher-order function *get* that returns the name and tasks of each element of a collection.

The query *Q* returns each department, the people of interest associated with that department, and their tasks. We assign the special task “buy” to clients.

```

Q = for (x ← organisation)
  return (⟨department = x.name,
         people = (get (x.employees) (λy.y.tasks))
              ∪ (get (x.contacts) (λy.“buy”))

```

(id)	dept
1	Product
2	Quality
3	Research
4	Sales

(id)	dept	employee	salary
1	Product	Alex	20000
2	Product	Bert	900
3	Research	Cora	50000
4	Research	Drew	60000
5	Sales	Erik	2000000
6	Sales	Fred	700
7	Sales	Gina	100000

(id)	employee	task
1	Alex	build
2	Bert	build
3	Cora	abstract
4	Cora	build
5	Cora	call
6	Cora	dissemble
7	Cora	enthuse
8	Drew	abstract
9	Drew	enthuse
10	Erik	call
11	Erik	enthuse
12	Fred	call
13	Gina	call
14	Gina	dissemble

(id)	dept	contact	client
1	Product	Pam	false
2	Product	Pat	true
3	Research	Rob	false
4	Research	Roy	false
5	Sales	Sam	false
6	Sales	Sid	false
7	Sales	Sue	true

Figure 1. Sample data

The result type of *Q* is:

```

Result = Bag ⟨department : String,
             people : Bag ⟨name : String, tasks : Bag String⟩

```

The result of running *Q* is:

```

[⟨department = “Product”,
  people = [⟨name = “Bert”, tasks = [“build”]⟩,
            ⟨name = “Pat”, tasks = [“buy”]⟩]⟩,
 ⟨department = “Research”, people = ∅⟩,
 ⟨department = “Sales”,
  people = [⟨name = “Erik”, tasks = [“call”, “enthuse”]⟩,
            ⟨name = “Fred”, tasks = [“call”]⟩,
            ⟨name = “Sue”, tasks = [“buy”]⟩]⟩]

```

How it really works The nested input data above can actually be stored in flat relational tables in a number of ways; one way is via the following schema Σ^b containing four tables:

```

Σb(tasks) = Bag ⟨employee : String, task : String⟩
Σb(employees) = Bag ⟨dept : String, employee : String,
                    salary : Int⟩
Σb(contacts) = Bag ⟨dept : String, contact : String,
                   client : Bool⟩
Σb(departments) = Bag ⟨dept : String⟩

```

To simplify the example, we assume in addition that every table has an integer-value *id* field, such that each row has a unique identifier; we omit it from the schema for readability. (In the rest of the paper, we will not require that each row has a unique identifier. Instead we provide a feature for generating unique indexes that can ultimately be translated into SQL’s row_number construct.) The tabular form of the nested data above is shown in Figure 1.

We can write a query Q_Σ that maps data in the flat schema Σ^b to the nested schema Σ :

```

for (x ← departments)
  return (⟨name = x.dept,
    employees =
      for (y ← employees) where (x.dept = y.dept)
        return (⟨name = y.name,
          tasks =
            for (z ← tasks)
              where (y.employee = z.employee)
                return z.task
          salary = y.salary⟩⟩)
    contacts =
      for (y ← contacts) where (x.dept = y.dept)
        return (⟨name = y.name, client = y.client⟩⟩)
  )

```

By composing Q with Q_Σ , we obtain a query of type $\Sigma^b \rightarrow Result$. According to previous work on query normalisation [18, 28], we know that any query that maps a *flat* input schema to a *flat* result type can be translated to a single SQL query that executes remotely on the database. Links, for example, already supports this [7, 13], and so if Q were a query that extracted a flat result from the nested schema Σ , we could compose with Q_Σ and normalise to obtain a flat query that can be executed on the database. Typically this is much more efficient than loading all of the data and running the query in-memory. We can normalise $Q \circ Q_\Sigma$ to obtain the following query Q' :

```

Q' = for (x ← departments)
  returna (⟨department = x.dept,
    people =
      (for (y ← employees) where (x.dept = y.dept ∧
        (y.salary < 1000 ∨ y.salary > 1000000))
        returnb (⟨name = y.name,
          tasks = for (z ← tasks)
            where (z.employee = y.employee)
              returnc z.task)
        )
      )
    )
  ⊔
  (for (y ← contacts) where (x.dept = y.dept ∧ y.client)
    returnd (⟨name = y.name,
      tasks = returne “buy”⟩)
  )

```

Notice that we have annotated each return with a unique tag $a, b, \dots \in Tag$. These annotations are explained shortly.

Now, however, we are faced with a problem: SQL databases do not directly support nested collections. However, in theory (see for example Van den Bussche [25]) one can translate a query $Q : \Sigma^b \rightarrow Result$ that maps flat input Σ^b to nested output $Result$ to a fixed number of flat queries $Q_1 : \Sigma^b \rightarrow Result_1^b, \dots, Q_n : \Sigma^b \rightarrow Result_n^b$ whose results can be combined via a mediating query $Q^\# : Result_1^b \times \dots \times Result_n^b \rightarrow Result$. Thus, we can simulate the query $Q \circ Q_\Sigma : \Sigma^b \rightarrow Result$ using $Q^\# \circ (Q_1, \dots, Q_n)$, where queries Q_1, \dots, Q_n can be run remotely on the database. The number of intermediate queries n is the *nesting degree* of $Result$, that is, the number of collection type constructors in the result type. For our example, the nesting degree of $Result$ is three, which means Q can be shredded into three flat queries.

The basic idea is straightforward. Whenever a nested bag appears in the output of a query, we generate an index that uniquely identifies the current context. Then a separate query produces the contents of the nested bag, where each element is paired up with its parent index. Each inner level of nesting requires a further query.

Let us now consider the top-level query Q_1 .

```

for (x ← departments)
  returna (⟨department = x.dept, people = ⟨a, x.id⟩⟩)

```

For each iteration, we store an index in place of the nested data. The index consists of a static component a (corresponding to the tag a on the comprehension) and a dynamic component id . The static component identifies which comprehension the index is associated with. The dynamic component identifies the current loop iteration. In this case a is redundant because there is only one top-level comprehension, but the need for static indexes will become apparent later. The result of running Q_1 is:

```

[⟨department = “Product”, people = ⟨a, 1⟩⟩,
 ⟨department = “Research”, people = ⟨a, 2⟩⟩,
 ⟨department = “Sales”, people = ⟨a, 3⟩⟩]

```

Now let us consider the auxiliary query Q_2 for generating the bag bound to the `people` field.

```

(for (x ← departments)
  for (y ← employees) where (x.dept = y.dept ∧
    (y.salary < 1000 ∨ y.salary > 1000000))
    returnb (⟨⟨a, x.id⟩,
      ⟨name = y.name, tasks = ⟨b, x.id, y.id⟩⟩))
  )
⊔
(for (x ← departments)
  for (y ← contacts) where (x.dept = y.dept ∧ y.client)
    returnd (⟨⟨a, x.id⟩,
      ⟨name = y.name, tasks = ⟨d, x.id, y.id⟩⟩))
  )

```

The idea here is that we can ultimately stitch the results of this query together with the results of Q_1 by joining the inner indexes of Q_1 to the outer indexes of Q_2 . In both cases the static components of these indexes are the same tag a .

More interestingly, this query also generates further indexes for the tasks associated with each person. The two halves of the union have different static indexes for the tasks b and d , because they arise from different comprehensions in the source term. Furthermore, the dynamic index now consists of two `id` fields ($x.id$ and $y.id$) in each half of the union.

The result of running Q_2 is:

```

[[⟨a, 1⟩, ⟨name = “Bert”, tasks = ⟨b, 1, 2⟩⟩],
 ⟨a, 3⟩, ⟨name = “Erik”, tasks = ⟨b, 3, 5⟩⟩],
 ⟨a, 3⟩, ⟨name = “Fred”, tasks = ⟨b, 3, 6⟩⟩],
 ⟨a, 1⟩, ⟨name = “Pat”, tasks = ⟨d, 1, 2⟩⟩],
 ⟨a, 3⟩, ⟨name = “Sue”, tasks = ⟨d, 3, 7⟩⟩]]

```

It is easy to see that joining the `people` field of Q_1 to the outer index of Q_2 correctly associates each person with the appropriate department.

Finally, let us consider the innermost query Q_3 for generating the bag bound to the `tasks` field.

```

(for (x ← departments)
  for (y ← employees) where (x.dept = y.dept ∧
    (y.salary < 1000 ∨ y.salary > 1000000))
    for (z ← tasks) where (z.employee = y.employee)
      returnc (⟨⟨b, x.id, y.id⟩, z.task⟩)
  )
⊔
(for (x ← departments)
  for (y ← contacts) where (x.dept = y.dept ∧ y.client)
    returne (⟨⟨d, x.id, y.id⟩, “buy”⟩)
  )

```

This query does not have any inner indexes because it is at the deepest level of nesting so does not construct any further nested collections. If it did, then we could use the static indexes c and e in the same way as before, but we would run into typing difficulties with the dynamic indexes, as the left branch iterates over three tables, whereas the right branch iterates over two. This limitation is

addressed by the more general indexing scheme used in Section 4. The result of running Q_3 is:

```
{(b, 1, 2), "build"}, (b, 3, 5), "call"}, (b, 3, 5), "enthuse"},
  (b, 3, 6), "call"}, (d, 1, 2), "buy"}, (d, 3, 7), "buy")}
```

Again, it is easy to see that joining the `tasks` field of Q_2 to the outer index of Q_3 correctly associates each task with the appropriate outlier.

Note that each of the queries Q_1, Q_2, Q_3 produces records that contain other records as fields. This is not strictly allowed by SQL, but it is straightforward to eliminate such nested records from a query with no nested collections in its result type, or re-introduce this nesting. Thus, we can convert each of Q_1, Q_2, Q_3 to queries that run remotely on any SQL database. For instance, Q_3 can be converted to the following SQL:

```
(select 2 as i1_1, x.id as i1_2, y.id as i1_3, z.task as i2
  from departments as x, employees as y, tasks as z
  where (x.employee = y.employee))
union all
(select 4 as i1_1, x.id as i1_2, y.id as i1_3, "buy" as i2
  from departments as x, contacts as y
  where (x.employee = y.employee))
```

Note that each static index a, b, c, \dots has been converted to a number 1, 2, 3, \dots , and records have been flattened by concatenating field names using `_` as a delimiter.

The three queries have the following result types:

```
Result1 = Bag ⟨department : String, people : Indexa⟩
Result2 = Bag ⟨Indexa, ⟨name : String, tasks : Indexb,d⟩⟩
Result3 = Bag ⟨Indexb,d, String⟩
```

where $Index_a = \langle Tag, Int \rangle$ and $Index_{b,d} = \langle Tag, Int, Int \rangle$. Once the results $R_1 : Result_1, R_2 : Result_2, R_3 : Result_3$ have been evaluated on the database, they are shipped back to the server where we can run the following code in-memory to *stitch* the three tables together into a single value: the result of the original nested query. The code for this follows the same idea as the query Q_Σ that constructs the nested version of Σ from Σ^b .

```
for (x ← Q1)
  return (⟨department = x.dept,
    people =
      for ((i, y) ← Q2) where (x.people = i))
    return (⟨name = y.name,
      tasks =
        for ((j, z) ← Q3) where (y.tasks = j)
          return z⟩))
```

We want to emphasise, however, that all of the code mentioned in this discussion is generated by the *implementation*, not by the programmer. All the programmer needs to do is write the concise, higher-order, nested query shown earlier.

We also want to emphasise that it is unnecessary to construct the nested version of the input data explicitly in-memory or on the database. Although the notional first step of our implementation approach is to apply query Q_Σ that constructs the nested database from the flat tables, there is no need to materialise this data structure anywhere. Instead, by composing with Q and normalising, we effectively deforest the intermediate nested data structure. In the common case where the query result is much smaller than the full database, this is usually much faster than shipping all of the data to the server.

In the rest of the paper we focus attention on the critical step of transforming a query Q whose input is flat and output is nested to multiple flat queries Q_1, \dots, Q_n , such that evaluating Q (in-memory) is equivalent to evaluating Q_1, \dots, Q_n on the database and then stitching the results together on the server.

1.2 Outline

The rest of the paper is organised as follows. Section 2 gives background and Section 3 describes query normalisation. Section 4 defines a translation from normal forms to shredded terms, *shredded packages* for bundling shredded queries or shredded results together, a semantics for shredded query packages, and a function for stitching shredded results back together to form a nested result. Section 5 shows that the shredding translation is correct. Section 6 gives a translation for providing flat indexes using let-insertion. Section 7 outlines a translation to SQL. Section 8 discusses variations and extensions. Section 9 discusses related work and Section 10 concludes.

The details of proofs are given in appendices. Also, for completeness, the details of the initial normalisation stage and the (standard) record-flattening stage are given in appendices.

2. Background

2.1 Notational conventions

We will use metavariables x, y, \dots, f, g for *variables*, and c, d, \dots for *constants* and primitive operations. We also use letters t, t', \dots for *table names*, $\ell, \ell', \ell_i, \dots$ for *record labels* and a, b, \dots for *tags*.

We write $M[x := N]$ for capture-avoiding substitution of N for x in M . We write \vec{x} for a vector x_1, \dots, x_n . Moreover, we extend vector notation pointwise to other constructs, writing, for example, $\langle \vec{\ell} = \vec{M} \rangle$ for a list of correlated pairs $\langle \ell_1 = M_1, \dots, \ell_n = M_n \rangle$.

We write: square brackets $[-]$ for the meta level list constructor; $w::\vec{v}$ for adding the element w onto the front of the list \vec{v} ; and $\vec{v}++\vec{w}$ for the result of appending the list \vec{w} onto the end of the list \vec{v} ; and *concat* for the function that concatenates a list of lists.

In the meta language we make extensive use of comprehensions, primarily list comprehensions, using a Haskell-like notation. For instance, $[v \mid x \leftarrow xs, y \leftarrow ys, p]$, returns a copy of v for each pair $\langle x, y \rangle$ of elements of xs and ys such that the predicate p holds. We write $[v_i]_{i=1}^n$ as shorthand for $[v \mid 1 \leq i \leq n]$ and similarly, e.g., $\langle \ell_i = M_i \rangle_{i=1}^n$ for $\langle \ell_1 = M_1, \dots, \ell_n = M_n \rangle$.

2.2 Nested relational calculus

We take the higher-order nested relational calculus (evaluated over bags) as our starting point. We call this λ_{NRC} . The types of λ_{NRC} include base types (integers, strings, booleans), record types $\langle \ell : \vec{A} \rangle$, bag types $\text{Bag } A$, and function types $A \rightarrow B$.

Types	$A, B ::= O \mid \langle \vec{\ell} : \vec{A} \rangle \mid \text{Bag } A \mid A \rightarrow B$
Base types	$O ::= \text{Int} \mid \text{Bool} \mid \text{String}$

We say that a type is *positive* if it contains no function types and *flat* if it is positive and contains no collection types.

The terms of λ_{NRC} include λ -abstractions, applications, and the standard terms of nested relational calculus.

Terms	$M, N ::= x \mid c(\vec{M}) \mid \text{table } t \mid \text{if } M \text{ then } N \text{ else } N'$
	$\mid \lambda x. M \mid M N \mid \langle \vec{\ell} = \vec{M} \rangle \mid M.\ell$
	$\mid \text{return } M \mid \emptyset \mid M \uplus N \mid \text{for } (x \leftarrow M) N$

We assume that the constants and primitive functions include boolean values with negation and conjunction, and integer values with standard arithmetic operations. We assume special labels $\#_1, \#_2, \dots$ and encode tuple types $\langle A_1, \dots, A_n \rangle$ as record types $\langle \#_1 : A_1, \dots, \#_n : A_n \rangle$, and similarly tuple terms $\langle M_1, \dots, M_n \rangle$ as record terms $\langle \#_1 = M_1, \dots, \#_n = M_n \rangle$. We assume fixed signatures $\Sigma(t)$ and $\Sigma(c)$ for tables and constants. The former are constrained to have flat relation type $(\text{Bag } \langle \ell_1 : O_1, \dots, \ell_n : O_n \rangle)$, and the latter to be first order n -ary functions $(\langle O_1, \dots, O_n \rangle \rightarrow O)$.

$$\begin{aligned}
\mathcal{N}[[x]]_\rho &= \rho(x) \\
\mathcal{N}[[c(X_1, \dots, X_n)]]_\rho &= [[c]](\mathcal{N}[[X_1]]_\rho, \dots, \mathcal{N}[[X_n]]_\rho) \\
\mathcal{N}[[\lambda x. M]]_\rho &= \lambda v. \mathcal{N}[[M]]_{\rho[x \mapsto v]} \\
\mathcal{N}[[M N]]_\rho &= \mathcal{N}[[M]]_\rho(\mathcal{N}[[N]]_\rho) \\
\mathcal{N}[[\langle \ell_i = M_i \rangle_{i=1}^n]]_\rho &= \langle \ell_i = \mathcal{N}[[M_i]]_\rho \rangle_{i=1}^n \\
\mathcal{N}[[M.\ell]]_\rho &= \mathcal{N}[[M]]_\rho.\ell \\
\mathcal{N}[[\text{if } L \text{ then } M \text{ else } N]]_\rho &= \begin{cases} \mathcal{N}[[M]]_\rho, & \text{if } \mathcal{N}[[L]]_\rho = \text{true} \\ \mathcal{N}[[N]]_\rho, & \text{if } \mathcal{N}[[L]]_\rho = \text{false} \end{cases} \\
\mathcal{N}[[\text{return } M]]_\rho &= [\mathcal{N}[[M]]_\rho] \\
\mathcal{N}[[\emptyset]]_\rho &= [] \\
\mathcal{N}[[M \uplus N]]_\rho &= \mathcal{N}[[M]]_\rho \uplus \mathcal{N}[[N]]_\rho \\
\mathcal{N}[[\text{for } (x \leftarrow M) N]]_\rho &= [\mathcal{N}[[N]]_{\rho[x \mapsto v]} \mid v \leftarrow \mathcal{N}[[M]]_\rho] \\
\mathcal{N}[[\text{table } t]]_\rho &= [[t]]
\end{aligned}$$

Figure 2. Semantics of higher-order nested relational queries

Most language constructs are standard. The \emptyset expression builds an empty bag, return M constructs a singleton, and $M \uplus N$ builds the bag union of two collections. The $\text{for } (x \leftarrow M) N$ comprehension construct iterates over a bag obtained by evaluating M , binds x to each element, evaluates N to another bag for each such binding, and takes the union of the results. (That is, it is a monadic bind or concat-map operation, not a map).

Semantics We give a set-theoretic denotational semantics in a Haskell-like meta-language. Though we give a bag semantics, for convenience we use lists as our primary meta-language collection type. Thus, we interpret object-level bags as meta-level lists. For meta-level values v and v' , we write $v \simeq v'$ if v and v' are equivalent up to permutation of list elements.

We interpret base types as base types, functions as functions, records as records, and bags as lists. For each table $t \in \text{dom}(\Sigma)$, we assume a fixed interpretation $[[t]]$ of t as a list of records of type $\Sigma(t)$. In SQL, tables do not have a list semantics by default, but we can impose one by choosing a canonical ordering for the rows of the table. The most general approach, which we adopt, is to order by all of the columns arranged in lexicographic order (assuming some ordering on field names). See Section 8 for a discussion of alternatives. The reason why we do need lists rather than bags here is that the order in which indexes are generated must be consistent across shredded queries.

We assume fixed interpretations $[[c]]$ for the constants and primitive operations. The (standard) typing rules and semantics of nested relational calculus are shown in Figure 2 and Figure 3; our previous paper [13] shows how to embed this query language into a general-purpose language using a type-and-effect system.

If $\Gamma \vdash M, N : A$ then we write $M \simeq N$ when $\mathcal{N}[[M]]_\rho \simeq \mathcal{N}[[N]]_\rho$ holds for all $\rho : \Gamma$. We let ρ range over environments mapping variables to values, writing ε for the empty environment and $\rho[x \mapsto v]$ for the extension of ρ with x bound to v .

3. Query normalisation

The first stage of our shredding algorithm is to rewrite the query to a normal form that has a number of useful properties. This is a simple adaptation of normalisation techniques in prior work [7, 13, 28]. Essentially, the idea is to apply η -expansion, β -reduction, and commuting conversions for bag comprehensions to reduce the query to a form with no λ -abstractions and no superfluous comprehensions.

For a *flat-flat* query (that is, one whose input and output types are flat tables with no nested collections), the normal forms pro-

duced by Cooper’s (and our) algorithm are:

Query terms	$L, M, N ::= \uplus \vec{C}$
Comprehensions	$C ::= \text{for } (\vec{G} \text{ where } X) \text{ return } R$
Generators	$G ::= x \leftarrow t$
Record terms	$R ::= \langle \ell = \vec{X} \rangle$
Base terms	$X ::= x.\ell \mid c(\vec{X})$

These normal forms correspond precisely¹ to a fragment of SQL:

Query terms	$L, M, N ::= (\text{union all}) \vec{C}$
Comprehensions	$C ::= \text{select } R \text{ from } \vec{G} \text{ where } X$
Generators	$G ::= t \text{ as } x$
Record terms	$R ::= \vec{X} \text{ as } \ell$
Base terms	$X ::= x.\ell \mid c(\vec{X})$

Thus, our normalisation algorithm provides an alternative proof of Cooper’s main result:

Theorem 1 (Cooper [7]). *Any flat-flat λ_{NRC} query can be translated to a single equivalent SQL query.*

In Links, query normalisation is an important part of the execution model. When a subexpression denoting a query is evaluated, the subexpression is first normalised and then converted to SQL, which is sent to the database for evaluation; the tuples received in response are then converted into a Links value and execution proceeds. Note that at run-time, the query expression is *closed* in the sense that all of its free variables have been replaced with values; however, it may still contain table names whose values are stored remotely in the database.

For *flat-nested* queries that read from flat tables and produce a nested result value, our normalisation procedure is similar, but in contrast to Cooper’s algorithm, we hoist all conditionals into the nearest enclosing comprehension as where clauses. The resulting normal forms are:

Query terms	$L ::= \uplus \vec{C}$
Comprehensions	$C ::= \text{for } (\vec{G} \text{ where } X) \text{ return } M$
Generators	$G ::= x \leftarrow t$
Normalised terms	$M, N ::= X \mid R \mid L$
Record terms	$R ::= \langle \ell = \vec{M} \rangle$
Base terms	$X ::= x.\ell \mid c(\vec{X})$

Any closed flat-nested query L can be converted to an equivalent term in the above normal form.

Theorem 2. *Given a closed flat-nested query $\vdash M : \text{Bag } A$, there exists a normalisation function $\text{norm}_{\text{Bag } A}$, such that $M \simeq \text{norm}_{\text{Bag } A}(M)$ is in normal form.*

The normalisation algorithm and correctness proof are similar to those in previous papers [7, 13], and are given in Appendix A. The normal forms for flat-nested queries do not correspond directly to SQL; the next two sections show how to bridge this gap.

4. Shredding and stitching

4.1 The shredding translation

In order to shred nested queries, we introduce an abstract type *Index* of *indexes* for maintaining the correspondence between outer and inner queries. An index $a \diamond d$ has a static component a and a dynamic component d . The static component identifies

¹Strictly speaking, this is not quite a subset of SQL as SQL cannot handle comprehensions in which R is a row with no columns, or empty unions. These idiosyncratic limitations are easy to circumvent.

$\frac{\text{VAR}}{\Gamma, x : A \vdash x : A}$	$\frac{\text{CONSTANT} \quad \Sigma(c) = \langle O_1, \dots, O_n \rangle \rightarrow O' \quad [\Gamma \vdash M_i : O_i]_{i=1}^n}{\Gamma \vdash c(M_1, \dots, M_n) : O'}$	$\frac{\text{LAM} \quad \Gamma, x : A \vdash M : B}{\Gamma \vdash \lambda x. M : A \rightarrow B}$
$\frac{\text{APP} \quad \Gamma \vdash M : A \rightarrow B \quad \Gamma \vdash N : A}{\Gamma \vdash M N : B}$	$\frac{\text{RECORD} \quad [\Gamma \vdash M_i : A_i]_{i=1}^n}{\Gamma \vdash \langle \ell_i = M_i \rangle_{i=1}^n : \langle \ell_i = A_i \rangle_{i=1}^n}$	$\frac{\text{PROJECT} \quad \Gamma \vdash M : \langle \ell_i : A_i \rangle_{i=1}^n}{\Gamma \vdash M.\ell_j : A_j}$
$\frac{\text{IF} \quad \Gamma \vdash M : \text{Bool} \quad \Gamma \vdash N : A \quad \Gamma \vdash N' : A}{\Gamma \vdash \text{if } M \text{ then } N \text{ else } N' : A}$	$\frac{\text{EMPTY}}{\Gamma \vdash \emptyset : \text{Bag } A}$	$\frac{\text{SINGLETON} \quad \Gamma \vdash M : A}{\Gamma \vdash \text{return } M : \text{Bag } A}$
$\frac{\text{UNION} \quad \Gamma \vdash M : \text{Bag } A \quad \Gamma \vdash N : \text{Bag } A}{\Gamma \vdash M \uplus N : \text{Bag } A}$	$\frac{\text{FOR} \quad \Gamma \vdash M : \text{Bag } A \quad \Gamma, x : A \vdash N : \text{Bag } B}{\Gamma \vdash \text{for } (x \leftarrow M) N : \text{Bag } B}$	$\frac{\text{TABLE} \quad \Sigma(t) = \text{Bag } \langle \vec{\ell} : \vec{O} \rangle}{\Gamma \vdash \text{table } t : \text{Bag } \langle \vec{\ell} : \vec{O} \rangle}$

Figure 3. Typing rules for higher-order nested queries

which comprehension the index is associated with. The dynamic component identifies the current loop iteration.

As a pre-processing step, we annotate each comprehension body in a normalised term with a unique name a — the static component of an index. We write the annotations as superscripts, for example:

$$\text{for } (\vec{G} \text{ where } X) \text{ return}^a M$$

Next, we modify types so that bag types have an explicit index component and we use indexes to replace nested occurrences of bags within other bags:

$$\begin{array}{ll} \text{Shredded types} & A, B ::= \text{Bag } \langle \text{Index}, F \rangle \\ \text{Flat types} & F ::= O \mid \langle \vec{\ell} : \vec{F} \rangle \mid \text{Index} \end{array}$$

We also adapt the syntax of terms to incorporate indexes. After shredding, terms will have the following forms:

$$\begin{array}{ll} \text{Query terms} & L, M ::= \uplus \vec{C} \\ \text{Comprehensions} & C ::= \text{return}^a \langle I, N \rangle \\ & \quad \mid \text{for } (\vec{G} \text{ where } X) C \\ \text{Generators} & G ::= x \leftarrow t \\ \text{Inner terms} & N ::= X \mid R \mid I \\ \text{Record terms} & R ::= \langle \vec{\ell} = \vec{N} \rangle \\ \text{Base terms} & X ::= x.\ell \mid c(\vec{X}) \\ \text{Indexes} & I, J ::= a \diamond d \\ \text{Dynamic indexes} & d ::= \text{out} \mid \text{in} \end{array}$$

A comprehension is now constructed from a sequence of generator clauses of the form $\text{for } (\vec{G} \text{ where } X)$ followed by a body of the form $\text{return}^a \langle I, N \rangle$. Each level of nesting gives rise to such a generator clause. The body always returns a pair $\langle I, N \rangle$ of an outer index I , denoting where the result values from the shredded query should be spliced into the final nested result, and a (flat) inner term N . Records are restricted to contain inner terms. Inner terms are either base types, records, or indexes, which replace nesting. We assume a distinguished top level static index \top , which allows us to treat all levels uniformly. Each shredded term is associated with an outer index out and an inner index in . In fact out always appears in the left component of a comprehension body, and only there, and in only appears in the right component of a comprehension body. These properties will become apparent when we specify the shredding transformation on terms.

An additional constraint that is not captured by the above grammar is that all comprehensions C_i in a shredded query $\uplus_{i=1}^n C_i$

must have the same non-zero *nesting degree*, that is, there exists $m > 0$ such that $\text{degree}(C_i) = m$ for $1 \leq i \leq n$, where:

$$\begin{aligned} \text{degree}(\text{return}^a \langle I, N \rangle) &= 0 \\ \text{degree}(\text{for } (\vec{G} \text{ where } X) C) &= 1 + \text{degree}(C) \end{aligned}$$

This constraint will be guaranteed by the shredding transformation.

To define the shredding transformation we use paths to point to parts of types. We use paths made up of the symbol \downarrow (representing traversing a bag constructor) and record labels ℓ .

$$\text{Paths } p ::= \epsilon \mid \downarrow.p \mid \ell.p$$

As paths are simply lists of \downarrow s and labels, we will sometimes write $p.\downarrow$ for the path p with \downarrow appended at the end and similarly for $p.\ell$. Furthermore, we will write $p.\vec{\ell}$ for the path p with all the labels of $\vec{\ell}$ appended. The function $\text{paths}(A)$ defines the set of paths to bags in a type A :

$$\begin{aligned} \text{paths}(O) &= \{\} \\ \text{paths}(\langle \ell_i : A_i \rangle_{i=1}^n) &= \{\ell_i.p \mid p \leftarrow \text{paths}(A_i)\}_{i=1}^n \\ \text{paths}(\text{Bag } A) &= \{\epsilon\} \cup \{\downarrow.p \mid p \leftarrow \text{paths}(A)\} \end{aligned}$$

The *nesting degree* of a type A is the number of paths to bag constructors in A , that is, $\text{degree}(A) = |\text{paths}(A)|$.

We now define a shredding translation on types. Given a path $p \in \text{paths}(A)$, the type $\llbracket A \rrbracket_p$ is the *outer shredding* of A at p , a shredded type that corresponds to the bag at path p in A . This is defined in terms of the *inner shredding* $\llbracket A \rrbracket$, a flat type that represents the contents of the bag.

$$\begin{aligned} \llbracket \text{Bag } A \rrbracket_\epsilon &= \text{Bag } \langle \text{Index}, \llbracket A \rrbracket \rangle \\ \llbracket \text{Bag } A \rrbracket_{\downarrow.p} &= \llbracket A \rrbracket_p \\ \llbracket \langle \vec{\ell} : \vec{A} \rangle \rrbracket_{\ell_i.p} &= \llbracket A_i \rrbracket_p \\ \llbracket O \rrbracket &= O \\ \llbracket \langle \ell_i : A_i \rangle_{i=1}^n \rrbracket &= \langle \ell_i : \llbracket A_i \rrbracket \rangle_{i=1}^n \\ \llbracket \text{Bag } A \rrbracket &= \text{Index} \end{aligned}$$

For example, consider the example result type *Result* from Section 1.1. Its nesting degree is 3, and its paths are:

$$\text{paths}(\text{Result}) = \{\epsilon, \downarrow.\text{people}.\epsilon, \downarrow.\text{people}.\downarrow.\text{tasks}.\epsilon\}$$

Moreover, we can shred it in three ways using these three paths:

$$\begin{aligned} \llbracket \text{Result} \rrbracket_\epsilon &= \langle \text{Index}, \langle \text{department} : \text{String}, \text{people} : \text{Index} \rangle \rangle \\ \llbracket \text{Result} \rrbracket_{\downarrow.\text{people}.\epsilon} &= \text{Bag } \langle \text{Index}, \langle \text{name} : \text{String}, \text{tasks} : \text{Index} \rangle \rangle \\ \llbracket \text{Result} \rrbracket_{\downarrow.\text{people}.\downarrow.\text{tasks}.\epsilon} &= \text{Bag } \langle \text{Index}, \text{String} \rangle \end{aligned}$$

$$\begin{aligned}
\llbracket L \rrbracket_p &= \biguplus (\llbracket L \rrbracket_{\top, p}^*) \\
\llbracket X \rrbracket_a &= X \\
\llbracket \langle \ell_i = M_i \rangle_{i=1}^n \rrbracket_a &= \langle \ell_i = \llbracket M_i \rrbracket_a \rangle_{i=1}^n \\
\llbracket L \rrbracket_a &= a \diamond \text{in} \\
\llbracket \biguplus_{i=1}^n C_i \rrbracket_{a, p}^* &= \text{concat}(\llbracket C_i \rrbracket_{a, p}^* \rrbracket_{i=1}^n) \\
\llbracket \langle \ell_i = M_i \rangle_{i=1}^n \rrbracket_{a, \ell_j \cdot p}^* &= \llbracket M_j \rrbracket_{a, p}^* \\
\llbracket \text{for } (\vec{G} \text{ where } X) \text{ return }^b M \rrbracket_{a, \epsilon}^* &= [\text{for } (\vec{G} \text{ where } X) \text{ return }^b \langle a \diamond \text{out}, \llbracket M \rrbracket_b \rangle] \\
\llbracket \text{for } (\vec{G} \text{ where } X) \text{ return }^b M \rrbracket_{a, \downarrow, p}^* &= [\text{for } (\vec{G} \text{ where } X) C \mid C \leftarrow \llbracket M \rrbracket_{b, p}^*]
\end{aligned}$$

Figure 4. Shredding translation on terms

The shredding translation on terms $\llbracket L \rrbracket_p$ is given in Figure 4. This takes a term L and a path p and gives a query $\llbracket L \rrbracket_p$ that computes a result of type $\llbracket A \rrbracket_p$, where A is the type of L . The auxiliary translation $\llbracket M \rrbracket_{a, p}^*$ returns the shredded comprehensions of M along path p with outer static index a . The auxiliary translation $\llbracket M \rrbracket_a$ produces a flat representation of M with inner static index a . Note that the shredding translation is linear in time and space.

As a basic sanity check, we show that well-formed normalised terms shred to well-formed shredded terms of the appropriate shredded types. We will show the full correctness of the shredding translation in Section 5. Typing rules for shredded terms are shown in Appendix B.

Theorem 3. *Suppose L is a normalised flat-nested query with $\vdash L : A$ and $p \in \text{paths}(A)$, then $\vdash \llbracket L \rrbracket_p : \llbracket A \rrbracket_p$.*

4.2 Shredded packages

To maintain the relationship between shredded terms and the structure of the nested result they are meant to construct, we use *shredded packages*. A shredded package \hat{A} is a nested type with annotations, denoted $(-)^{\alpha}$, attached to each bag constructor.

$$\hat{A} ::= O \mid \langle \ell : \hat{A} \rangle \mid (\text{Bag } \hat{A})^{\alpha}$$

For a given package, the annotations are drawn from the same set. We write $\hat{A}(S)$ to denote a shredded package with annotations drawn from the set S . We write Λ_S for the set of shredded terms, and \mathcal{T}_S for the set of shredded types. We sometimes omit the type parameter when it is clear from context. Typing rules for shredded packages are shown in Figure 5.

Given a shredded package \hat{A} , we can erase its annotations to obtain its underlying type.

$$\begin{aligned}
\text{erase}(O) &= O \\
\text{erase}(\langle \ell_i : \hat{A}_i \rangle_{i=1}^n) &= \langle \ell_i : \text{erase}(\hat{A}_i) \rangle_{i=1}^n \\
\text{erase}((\text{Bag } \hat{A})^{\alpha}) &= \text{Bag } \text{erase}(\hat{A})
\end{aligned}$$

Given a shredded package $\hat{A}(S)$ and a function $f : S \rightarrow T$, we can map f over the annotations to obtain a new shredded package $\hat{A}'(T)$ such that $\text{erase}(\hat{A}) = \text{erase}(\hat{A}')$.

$$\begin{aligned}
\text{pmap}_f(O) &= O \\
\text{pmap}_f(\langle \ell_i : \hat{A}_i \rangle_{i=1}^n) &= \langle \ell_i : \text{pmap}_f(\hat{A}_i) \rangle_{i=1}^n \\
\text{pmap}_f((\text{Bag } \hat{A})^{\alpha}) &= (\text{Bag } \text{pmap}_f(\hat{A}))^{f(\alpha)}
\end{aligned}$$

Given a type A and a shredding function $f : \text{paths}(A) \rightarrow S$, we can construct a shredded package $\hat{A}(S)$.

$$\begin{aligned}
\text{package}_f(A) &= \text{package}_{f, \epsilon}(A) \\
\text{package}_{f, p}(O) &= O \\
\text{package}_{f, p}(\langle \ell_i : \hat{A}_i \rangle_{i=1}^n) &= \langle \ell_i : \text{package}_{f, p, \ell_i}(\hat{A}_i) \rangle_{i=1}^n \\
\text{package}_{f, p}(\text{Bag } A) &= (\text{Bag } \text{package}_{f, p, \downarrow}(A))^{f(p)}
\end{aligned}$$

Using *package*, we lift the type-level and term-level shredding functions $\llbracket - \rrbracket_p$ to produce shredded packages, where each annotation contains the shredded version of the input type or query along

$$\begin{array}{c}
\text{BASE} \\
\hline
\vdash O(\Lambda_S) : O(\mathcal{T}_S) \\
\\
\text{RECORD} \\
\hline
\frac{[- \hat{A}_i(\Lambda_S) : \hat{A}'_i(\mathcal{T}_S)]_{i=1}^n}{\vdash \langle \ell_i : \hat{A}_i(\Lambda_S) \rangle_{i=1}^n : \langle \ell_i : \hat{A}'_i(\mathcal{T}_S) \rangle_{i=1}^n} \\
\\
\text{BAG} \\
\hline
\frac{\vdash \hat{A}(\Lambda_S) : \hat{A}'(\mathcal{T}_S) \quad \vdash L : A \quad \text{erase}(\hat{A}(\Lambda_S)) = A}{\vdash (\text{Bag } \hat{A}(\Lambda_S))^L : (\text{Bag } \hat{A}'(\mathcal{T}_S))^A}
\end{array}$$

Figure 5. Typing rules for shredded packages

the path to the associated bag constructor.

$$\begin{aligned}
\text{shred}_B(A) &= \text{package}_{(\llbracket B \rrbracket_p)}(A) \\
\text{shred}_L(A) &= \text{package}_{(\llbracket L \rrbracket_p)}(A)
\end{aligned}$$

Let us now illustrate the shredding process by reference to the example query from the introduction. Shredding the *Result* type gives:

$$\begin{aligned}
\text{shred}_{\text{Result}}(\text{Result}) &= \\
&\text{Bag } \langle \text{department} : \text{String}, \\
&\quad \text{people} : \text{Bag } \langle \text{name} : \text{String}, \\
&\quad \quad \text{tasks} : (\text{Bag } \text{String})^{A_3} \rangle^{A_2} \rangle^{A_1}
\end{aligned}$$

where:

$$\begin{aligned}
A_1 &= \text{Bag } \langle \text{Index}, \langle \text{department} : \text{String}, \text{people} : \text{Index} \rangle \rangle \\
A_2 &= \text{Bag } \langle \text{Index}, \langle \text{name} : \text{String}, \text{tasks} : \text{Index} \rangle \rangle \\
A_3 &= \text{Bag } \langle \text{Index}, \text{String} \rangle
\end{aligned}$$

Shredding the query Q' gives the same package, except the type annotations A_1, A_2, A_3 become queries Q_1, Q_2, Q_3 . We will use Q_2 as a running example for the rest of the paper.

$$\begin{aligned}
Q_2 &= (\text{for } (x \leftarrow \text{departments}) \\
&\quad \text{for } (y \leftarrow \text{employees}) \text{ where } (x.\text{dept} = y.\text{dept} \wedge \\
&\quad \quad (y.\text{salary} < 1000 \vee y.\text{salary} > 1000000)) \\
&\quad \text{return}^b (\langle a \diamond \text{out}, (\text{name} = y.\text{name}, \text{tasks} = b \diamond \text{in}) \rangle)) \\
&\quad \biguplus (\text{for } (x \leftarrow \text{departments}) \\
&\quad \quad \text{for } (y \leftarrow \text{contacts}) \text{ where } (x.\text{dept} = y.\text{dept} \wedge y.\text{client}) \\
&\quad \quad \text{return}^d (\langle a \diamond \text{out}, (\text{name} = y.\text{name}, \text{tasks} = d \diamond \text{in}) \rangle))
\end{aligned}$$

Again, as a sanity check we show that erasure is the left inverse of type shredding and that term-level shredding operations preserve types.

Theorem 4. *Given a nested type A , we have $\text{erase}(\text{shred}_A(A)) = A$. Furthermore, if L is a closed, normalised, flat-nested query such that $\vdash L : A$ then $\vdash \text{shred}_L(A) : \text{shred}_A(A)$ also.*

5. Correctness

The main result we wish to prove is that if we shred an annotated normalised nested query, run all the resulting shredded queries, and stitch the shredded results back together, then that is the same as running the nested query directly.

The proof idea is as follows: a) define shredding on nested values; b) show that shredding commutes with query execution;

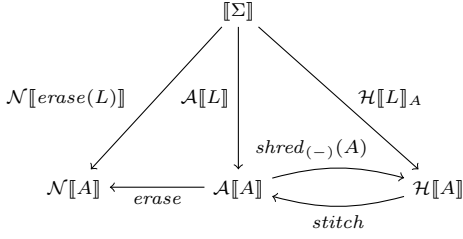


Figure 6. Correctness of shredding and stitching

and c) show that stitching is a left-inverse to shredding of values. Unfortunately the approach fails, because nested values do not contain enough information about the structure of the source query in order to generate the same indexes that are generated by shredded queries.

To fix the problem, we will give an annotated semantics $\mathcal{A}[-]$ that is defined only on queries that have first been converted to normal form and annotated with static indexes as described in Section 4.1. We will ensure that $\mathcal{N}[\text{erase}(L)] = \text{erase}(\mathcal{A}[L])$, where we overload the *erase* function to erase annotations on terms and values. The structure of the resulting proof is depicted in the diagram in Figure 6, where we view a query as a function from the interpretation of the Σ to the interpretation of its result type. To prove correctness is to prove that this diagram commutes.

5.1 Annotated semantics of nested queries

We annotate each bag element in the result set with an index.

Results	$s ::= [v_1 @ I_1, \dots, v_m @ I_m]$
Inner values	$w ::= c \mid r \mid I$
Rows	$r ::= \langle \ell_1 = w_1, \dots, \ell_n = w_n \rangle$
Indexes	I, J
Dynamic indexes	$\iota ::= 1 \mid \iota.i$

Indexes themselves are still abstract, but we will construct them from a concrete representation of static and dynamic indexes. Concretely, we represent each dynamic index ι as a list of integers. Each integer stores the current loop iteration around a comprehension. The advantage of this particular representation is that it admits a simple declarative semantics, which will make it straightforward to relate the nested semantics with a suitable semantics for shredded queries. The top level dynamic index is always 1.

The *canonical* choice for representing indexes is to take $I = a \diamond \iota$. We also admit alternative indexing schemes, by parameterising the semantics over an indexing function *index* mapping canonical indexes a concrete representation. In the simplest case, we take *index* to be the identity function. More generally, it may depend on the particular query being shredded. Informally, the only constraints on *index* are that it is defined on every canonical index $a \diamond \iota$ that we might need in order to shred a query, and it is injective on all canonical indexes. We will formalise these constraints in Section 5.5. For now we assume canonical indexes.

Given an indexing function *index*, the semantics of nested queries on annotated values is defined as follows:

$$\begin{aligned} \mathcal{A}[L] &= \mathcal{A}[L]_{\varepsilon, 1} \\ \mathcal{A}[\langle \bigoplus_{i=1}^n C_i \rangle_{\rho, \iota}] &= \text{concat}([\mathcal{A}[C_i]_{\rho, \iota}]_{i=1}^n) \\ \mathcal{A}[\langle \ell_i = M_i \rangle_{i=1}^n]_{\rho, \iota} &= \langle \ell_i = \mathcal{A}[M_i]_{\rho, \iota} \rangle_{i=1}^n \\ \mathcal{A}[X]_{\rho, \iota} &= \mathcal{N}[X]_{\rho} \\ \mathcal{A}[\text{for } ([x_i \leftarrow t_i]_{i=1}^n \text{ where } X) \text{ return }^a M]_{\rho, \iota} &= \\ &[\mathcal{A}[M]_{\rho[x_i \mapsto r_i]_{i=1, \dots, j}} @ \text{index}(a \diamond \iota, j) \\ & \mid \langle j, \bar{r} \rangle \leftarrow \text{enum}([\bar{r} \mid [r_i \leftarrow [t_i]_{i=1}^n, \mathcal{N}[X]_{\rho[x_i \mapsto r_i]_{i=1}^n}])]] \end{aligned}$$

As well as an environment, the current dynamic index is threaded through the semantics.

The *enum* function takes a list of elements and returns the same list with the element number paired up with each source element.

$$\text{enum}([v_1, \dots, v_m]) = [\langle 1, v_1 \rangle, \dots, \langle m, v_m \rangle]$$

The index annotations $@I$ on collection elements are needed solely for our correctness proof, and do not need to be materialised at run time.

5.2 Shredding and stitching annotated values

To define the semantics of shredded queries and packages, we use annotated values in which collections are annotated pairs of indexes and annotated values. Again, we allow annotations $@J$ on elements of collections to facilitate the proof of correctness. Typing rules for these values are shown in Appendix B.

$$\text{Results } s ::= [\langle I_1, w_1 \rangle @ J_1, \dots, \langle I_m, w_m \rangle @ J_m]$$

Shredding values Having defined suitably annotated versions of nested and shredded values, we now overload the shredding function to operate on nested values.

$$\begin{aligned} \llbracket s \rrbracket_p &= \llbracket s \rrbracket_{\top \diamond 1, p}^* \\ \llbracket [v_i @ J_i]_{i=1}^n \rrbracket_{I, \varepsilon}^* &= [\langle I, \llbracket v_i \rrbracket_{J_i}^* \rangle @ J_i]_{i=1}^n \\ \llbracket [v_i @ J_i]_{i=1}^n \rrbracket_{I, \downarrow, p}^* &= \text{concat}([\llbracket v_i \rrbracket_{J_i, p}^*]_{i=1}^n) \\ \llbracket \langle \ell_i = v_i \rangle_{i=1}^n \rrbracket_{I, \ell_i, p}^* &= \llbracket v_i \rrbracket_{I, p}^* \\ \llbracket c \rrbracket_I &= c \\ \llbracket \langle \ell_i = v_i \rangle_{i=1}^n \rrbracket_I &= \langle \ell_i = \llbracket v_i \rrbracket_I \rangle_{i=1}^n \\ \llbracket s \rrbracket_I &= I \end{aligned}$$

Note that in the cases for bags, the index annotation J is passed as an argument to the recursive call: this is where we need the ghost indexes, to relate the semantics of nested and shredded queries.

We lift the nested result shredding function to build an (annotated) shredded value package in the same way that we did for nested types and nested queries.

$$\begin{aligned} \text{shred}_s(\text{Bag } A) &= \text{package}_{(\llbracket s \rrbracket_-)}(\text{Bag } A) \\ \text{shred}_{s, p}(\text{Bag } A) &= \text{package}_{(\llbracket s \rrbracket_{-, p})}(\text{Bag } A) \end{aligned}$$

Stitching values A shredded value package can be stitched back together into a nested value as follows.

$$\begin{aligned} \text{stitch}(\hat{A}) &= \text{stitch}_{\top \diamond 1}(\hat{A}) \\ \text{stitch}_c(O) &= c \\ \text{stitch}_{\langle \ell_i = w_i \rangle_{i=1}^n}(\langle \ell_i : \hat{A}_i \rangle_{i=1}^n) &= \langle \ell_i = \text{stitch}_{w_i}(\hat{A}_i) \rangle_{i=1}^n \\ \text{stitch}_I((\text{Bag } \hat{A})^s) &= [(\text{stitch}_w(\hat{A})) @ J \\ & \mid \langle I', w \rangle @ J \leftarrow s, I' = I] \end{aligned}$$

The inner value parameter w to the auxiliary function $\text{stitch}_w(-)$ specifies which values to stitch along the current path.

5.3 Annotated semantics of shredded queries and packages

We now show how to interpret shredded queries and query packages over shredded values. The semantics of shredded queries is given in Figure 7. The semantics of a shredded query package is a *shredded value package* containing indexed results for each shredded query.

$$\begin{aligned} \mathcal{H}[A] &= \text{shred}_A(A) \\ \text{If } \vdash L : A, \text{ then: } & \mathcal{H}[L]_A : \mathcal{H}[A] \\ & \mathcal{H}[L]_A = \text{pmap}_{S[-]} \text{shred}_L(A) \end{aligned}$$

$$\begin{aligned}
\mathcal{S}[[L]] &= \mathcal{S}[[L]]_{\varepsilon,1} & \mathcal{S}[\langle \ell = N \rangle_{i=1}^n]_{\rho,\iota} &= \langle \ell_i = \mathcal{S}[[N_i]]_{\rho,\iota} \rangle_{i=1}^n & \mathcal{S}[[a \diamond \text{out}]]_{\rho,\iota,i} &= \text{index}(a \diamond \iota) \\
& & \mathcal{S}[[X]]_{\rho,\iota} &= \mathcal{N}[[X]]_{\rho} & \mathcal{S}[[a \diamond \text{in}]]_{\rho,\iota,i} &= \text{index}(a \diamond \iota.i) \\
\mathcal{S}[[\biguplus_{i=1}^n C_i]]_{\rho,\iota} &= \text{concat}([\mathcal{S}[[C_i]]_{\rho,\iota}]_{i=1}^n) & & & \mathcal{S}[\text{return}^a N]_{\rho,\iota} &= [\mathcal{S}[[N]]_{\rho,\iota} @ \text{index}(a \diamond \iota)] \\
\mathcal{S}[[\text{for}([x_i \leftarrow t_i]_{i=1}^n \text{ where } X) C]]_{\rho,\iota} &= \text{concat}([\mathcal{S}[[C]]_{\rho[x_i \mapsto r_i]_{i=1}^n, \iota, j} \mid \langle j, \vec{r} \rangle \leftarrow \text{enum}([\vec{r} \mid [r_i \leftarrow \llbracket t_i \rrbracket]_{i=1}^n, \mathcal{N}[[X]]_{\rho[x_i \mapsto r_i]_{i=1}^n}])])
\end{aligned}$$

Figure 7. Semantics of shredded queries

For example, here is the shredded package that we obtain after running the query Q' :

$$\mathcal{H}[[Q']]_A = \text{Bag} \langle \text{department} : \text{String}, \\
\text{people} : \text{Bag} \langle \text{name} : \text{String}, \\
\text{tasks} : (\text{Bag } \text{String})^{s_3} \rangle^{s_2} \rangle^{s_1}$$

where:

$$\begin{aligned}
s_1 &= [\langle \top \diamond 1, \langle \text{department} = \text{"Product"}, \text{people} = a \diamond 1.1 \rangle \rangle, \\
&\quad \langle \top \diamond 1, \langle \text{department} = \text{"Research"}, \text{people} = a \diamond 1.2 \rangle \rangle, \\
&\quad \langle \top \diamond 1, \langle \text{department} = \text{"Sales"}, \text{people} = a \diamond 1.3 \rangle \rangle] \\
s_2 &= [\langle a \diamond 1.1, \langle \text{name} = \text{"Bert"}, \text{tasks} = b \diamond 1.1.2 \rangle \rangle, \\
&\quad \langle a \diamond 1.3, \langle \text{name} = \text{"Erik"}, \text{tasks} = b \diamond 1.3.5 \rangle \rangle, \dots] \\
s_3 &= [\langle b \diamond 1.1.2, \text{"build"} \rangle, \langle b \diamond 1.3.5, \text{"call"} \rangle, \dots]
\end{aligned}$$

5.4 Main results

Due to space limits, we relegate the full details of the proof to Appendix C. There are three key theorems. The first states that shredding commutes with the annotated semantics.

Theorem 5. *If $\vdash L : \text{Bag } A$ then $\mathcal{H}[[L]]_{\text{Bag } A} = \text{shred}_{\mathcal{A}[[L]]}(\text{Bag } A)$.*

In order to allow shredded results to be correctly stitched together, we need the indexes at the end of each path to a bag through a nested value to be unique. We define $\text{indexes}_p(v)$, the indexes of nested value v along path p as follows.

$$\begin{aligned}
\text{indexes}_\varepsilon([v_i @ J_i]_{i=1}^n) &= [J_i]_{i=1}^n \\
\text{indexes}_{\downarrow.p}([v_i @ J_i]_{i=1}^n) &= \text{concat}([\text{indexes}_p(v_i)]_{i=1}^n) \\
\text{indexes}_{\ell_i.p}(\langle \ell_i = v_i \rangle_{i=1}^n) &= \text{indexes}_p(v_i)
\end{aligned}$$

We say that a nested value v is *well-indexed* if there exists a type A such that $\vdash v : A$, and for every path p in $\text{paths}(A)$, the elements of $\text{indexes}_p(v)$ are distinct.

Lemma 6. *If $\vdash L : A$, then $\mathcal{A}[[L]]$ is well-indexed.*

Our next theorem states that for well-indexed values, stitching is a left-inverse of shredding.

Theorem 7. *If $\vdash s : \text{Bag } A$ and s is well-indexed then $\text{stitch}(\text{shred}_s(\text{Bag } A)) = s$.*

Combining Theorem 5, Lemma 6 and Theorem 7 we obtain the main correctness result (see Figure 6).

Theorem 8. *If $\vdash L : \text{Bag } A$ then: $\text{stitch}(\mathcal{H}[[L]]_{\text{Bag } A}) = \mathcal{A}[[L]]$, and in particular: $\text{erase}(\text{stitch}(\mathcal{H}[[L]]_{\text{Bag } A})) = \mathcal{N}[[L]]$.*

5.5 Alternative indexing schemes

In order to formalise *valid* indexing schemes, we first define a function for computing the canonical indexes of a nested query.

$$\mathcal{I}[[L]] = \mathcal{I}[[L]]_{\varepsilon,1}$$

$$\begin{aligned}
\mathcal{I}[[\biguplus_{i=1}^n C_i]]_{\rho,\iota} &= \text{concat}([\mathcal{I}[[C_i]]_{\rho,\iota}]_{i=1}^n) \\
\mathcal{I}[\langle \ell_i = M_i \rangle_{i=1}^n]_{\rho,\iota} &= \text{concat}([\mathcal{I}[[M_i]]_{\rho,\iota}]_{i=1}^n) \\
\mathcal{I}[[X]]_{\rho,\iota} &= []
\end{aligned}$$

$$\begin{aligned}
\mathcal{I}[[\text{for}([x_i \leftarrow t_i]_{i=1}^n \text{ where } X) \text{return}^a M]]_{\rho,\iota} &= \\
&\text{concat}([\mathcal{I}[[M]]_{\rho[x_i \mapsto r_i]_{i=1}^n, \iota, j} \mid \langle j, \vec{r} \rangle \leftarrow \text{enum}([\vec{r} \mid [r_i \leftarrow \llbracket t_i \rrbracket]_{i=1}^n, \mathcal{N}[[X]]_{\rho[x_i \mapsto r_i]_{i=1}^n}])])
\end{aligned}$$

Note that $\mathcal{I}[-]$ follows essentially the same form as $\mathcal{A}[-]$, but instead of the nested value v it computes all indexes of v .

We define alternative indexing schemes by instantiating the *index* parameter of the shredded semantics (see Section 5.3). An *index* parameter is *valid* with respect to the closed nested query L if it is injective and defined on every canonical index in $\mathcal{I}[[L]]$.

Lemma 9. *If index is valid for L then $\mathcal{A}[[L]]$ is well-indexed.*

The only requirement on indexes in the proof of Theorem 8 is that nested values be well-indexed, hence the proof extends to any valid indexing scheme. We briefly mention two alternative valid indexing schemes: *flat indexes* and *natural indexes*.

Flat indexes The idea of flat indexes is to enumerate all of the canonical dynamic indexes associated with each static index and use the enumeration as the dynamic index.

Let ι be the i -th element of the list $[l' \mid a \diamond l' \leftarrow \mathcal{I}[[L]]]$, then $\text{index}_L^b(a \diamond \iota) = \langle a, i \rangle$. The flat indexing scheme is defined by setting $\text{index} = \text{index}_L^b$. Let $\mathcal{I}^b[[L]] = [\text{index}_L^b(I) \mid I \leftarrow \mathcal{I}[[L]]]$.

We elaborate on flat indexes in the next two sections, illustrating how they translate to SQL. This depends on the `row_number` operator, which we use to implement the functionality of *enum*.

Natural indexes Natural indexes are synthesised from row data. In order to generate a natural index for a query every table must have a primary key, that is, a collection of fields guaranteed to be unique for every row in the table. The type of natural indexes is a sum or variant type. The static index is a tag, and its payload is a tuple of primary keys.

Given a table t , let key_t be the function that given a row r of t returns the primary key of r . We now define a function to compute the list of natural indexes for a query L .

$$\mathcal{I}^b[[L]] = \mathcal{I}^b[[L]]_{\varepsilon,1}$$

$$\begin{aligned}
\mathcal{I}^b[[\biguplus_{i=1}^n C_i]]_{\rho,\iota} &= \text{concat}([\mathcal{I}^b[[C_i]]_{\rho,\iota}]_{i=1}^n) \\
\mathcal{I}^b[\langle \ell_i = M_i \rangle_{i=1}^n]_{\rho,\iota} &= \text{concat}([\mathcal{I}^b[[M_i]]_{\rho,\iota}]_{i=1}^n) \\
\mathcal{I}^b[[X]]_{\rho,\iota} &= []
\end{aligned}$$

$$\begin{aligned}
\mathcal{I}^b[[\text{for}([x_i \leftarrow t_i]_{i=1}^n \text{ where } X) \text{return}^a M]]_{\rho,\iota} &= \\
&\text{concat}([\mathcal{I}^b[[M]]_{\rho[x_i \mapsto r_i]_{i=1}^n, \iota, j} \mid \langle j, \vec{r} \rangle \leftarrow \text{enum}([\vec{r} \mid [r_i \leftarrow \llbracket t_i \rrbracket]_{i=1}^n, \mathcal{N}[[X]]_{\rho[x_i \mapsto r_i]_{i=1}^n}])])
\end{aligned}$$

If $a \diamond \iota$ is the i -th element of $\mathcal{I}[[L]]$, then $\text{index}_L^b(a \diamond \iota)$ is defined as the i -th element of $\mathcal{I}^b[[L]]$. The natural indexing scheme is defined by setting $\text{index} = \text{index}_L^b$.

An advantage of natural indexes is that when translating to SQL it is unnecessary to use the `row_number` function. Consequently, for a given comprehension all where clauses can be amalgamated (using the \wedge operator) and no auxiliary subqueries are needed. The downside is that the type of a dynamic index, which is determined by the tables being queried, may vary across the component comprehensions of a shredded query. Variant types can be used to fix this problem.

There are two ways of encoding variants in SQL. The first encoding packs variants into a record. The record consists of a distinguished tag column, and one column for each column of

each shredded variant argument in the variant type. The tag column holds an integer encoding which variant tag is active and the columns associated with the encoded tag are populated with their values. The other columns can contain nulls or any other values. The second encoding generates multiple queries: one for each tag.

We choose to focus on flat indexes in the remainder of the paper because they are more generally applicable than natural indexes, which require all rows of the source tables to be distinct.

6. Flat indexes and let-insertion

Our semantics for shredded queries uses canonical indexes. We now specify an object language that builds in flat indexes and move closer towards an SQL implementation. In order to do so, we introduce let-bound sub-queries, and translate each comprehension into the following form:

$$\text{let } q = \text{for } (\overrightarrow{G_{\text{out}}} \text{ where } X_{\text{out}}) \text{ return } N_{\text{out}} \text{ in} \\ \text{for } (\overrightarrow{G_{\text{in}}} \text{ where } X_{\text{in}}) \text{ return } N_{\text{in}}$$

The special index expression is available in each loop body, and is bound to the current iteration.

Following let-insertion, the types are as before, except indexes are represented as pairs of integers.

$$\begin{array}{ll} \text{Types} & A, B ::= \text{Bag } \langle \langle \text{Int}, \text{Int} \rangle, F \rangle \\ \text{Flat types} & F ::= O \mid \langle \ell : \vec{F} \rangle \mid \langle \text{Int}, \text{Int} \rangle \end{array}$$

Static indexes are represented as plain integers and dynamic integers as n -tuples. An index is a pair of a static index and a dynamic index.

The syntax of terms is adapted as follows:

$$\begin{array}{ll} \text{Query terms} & L, M ::= \uplus \vec{C} \\ \text{Comprehensions} & C ::= \text{let } q = S \text{ in } S' \\ \text{Subqueries} & S ::= \text{for } (\vec{G} \text{ where } X) \text{ return } N \\ \text{Data sources} & u ::= t \mid q \\ \text{Generators} & G ::= x \leftarrow u \\ \text{Inner terms} & N ::= X \mid R \mid \text{index} \\ \text{Record terms} & R ::= \langle \vec{\ell} = \vec{N} \rangle \\ \text{Base terms} & X ::= x.\ell \mid c(\vec{X}) \end{array}$$

Without loss of generality we α -rename all the variables in our source query to ensure that all bound variables have distinct names, and that none coincides with the distinguished name z used for let-bindings. The let-insertion translation \mathbf{L} is defined in Figure 8. Each comprehension is squashed into two sub-queries. The first generates the outer index. The second computes the results.

For example, applying \mathbf{L} to Q_2 from Section 4.2 yields:

$$\text{(let } q = \text{for } (x \leftarrow \text{departments}) \text{ return } \langle \langle \text{dept} = x.\text{dept} \rangle, \text{index} \rangle \text{ in} \\ \text{for } (z \leftarrow q, y \leftarrow \text{employees}) \text{ where } (z.1.\text{dept} = y.\text{dept} \wedge \\ (y.\text{salary} < 1000 \vee y.\text{salary} > 1000000)) \\ \text{return } (\langle \langle a, z.2 \rangle, \langle \text{name} = y.\text{name}, \text{tasks} = \langle b, \text{index} \rangle \rangle \rangle)) \\ \uplus$$

$$\text{(let } q = \text{for } (x \leftarrow \text{departments}) \text{ return } \langle \langle \text{dept} = x.\text{dept} \rangle, \text{index} \rangle \text{ in} \\ \text{for } (z \leftarrow q, y \leftarrow \text{contacts}) \text{ where } (z.1.\text{dept} = y.\text{dept} \wedge y.\text{client}) \\ \text{return } (\langle \langle a, z.2 \rangle, \langle \text{name} = y.\text{name}, \text{tasks} = \langle d, \text{index} \rangle \rangle \rangle))$$

To state and prove the correctness of let-insertion, we adapt the language of values as follows to incorporate concrete indexes:

$$\begin{array}{ll} \text{Results} & s ::= \vec{r} \\ \text{Inner values} & w ::= c \mid r \\ \text{Rows} & r ::= \langle \ell_1 = w_1, \dots, \ell_n = w_n \rangle \end{array}$$

The semantics is given in Figure 9. Rather than maintaining a canonical index, it generates a flat index for each subquery.

As a sanity check, we show that the translation is type-preserving:

Theorem 10. *If L is a shredded query such that $\vdash L : \text{Bag } \langle \text{Index}, F \rangle$, then $\vdash \mathbf{L}(L) : \text{Bag } \langle \langle \text{Int}, \text{Int} \rangle, F' \rangle$ (where F' is the result of replacing Index with $\langle \text{Int}, \text{Int} \rangle$ in F).*

In order to validate let-insertion, we need to show that the shredded semantics and let-inserted semantics agree.

Theorem 11. *Let $\vdash L : A$ and $\llbracket L \rrbracket_p = M$, and set $\text{index} = \text{index}_L^b$, then $\text{erase}(\mathcal{S}\llbracket M \rrbracket) = \mathcal{L}\llbracket \mathbf{L}(M) \rrbracket$.*

Proof sketch. The high-level idea is to separate results into data and indexes and compare each separately. It is straightforward, albeit tedious, to show that the definitions collapse to the same thing if we replace all dynamic indexes by unit. It then remains to show that the dynamic indexes agree. The pertinent case is the translation of a comprehension:

$$[\text{for } (\vec{G}_i \leftarrow X_i)]_{i=1}^n \text{ for } (\vec{G}_{\text{in}} \leftarrow X_{\text{in}}) \text{ return }^b \langle a \diamond \text{out}, N \rangle$$

which becomes:

$$\text{let } q = S_{\text{out}} \text{ in } S_{\text{in}}$$

for suitable S_{out} and S_{in} . The dynamic indexes computed by S_{out} coincide exactly with those of $\mathcal{I}^b\llbracket L \rrbracket$ at static index a , and the dynamic indexes computed by S_{in} , if there are any, coincide exactly with those of $\mathcal{I}^b\llbracket L \rrbracket$ at static index b . \square

7. Conversion to SQL

SQL does not support nested records. Our final translation flattens records. (In fact, we could have performed flattening at any earlier point, but leave it until the end for convenience so that we can use nested records in the earlier phases.) The (standard) details are presented in Appendix D for completeness.

In order to interpret shredded, flattened, let-inserted terms as SQL, we interpret index generators using SQL's OLAP facilities.

$$\begin{array}{ll} \text{Query terms} & L ::= (\text{union all}) \vec{C} \\ \text{Comprehensions} & C ::= \text{with } q \text{ as } (S) C \mid S' \\ \text{Subqueries} & S ::= \text{select } R \text{ from } \vec{G} \text{ where } X \\ \text{Data sources} & u ::= t \mid q \\ \text{Generators} & G ::= u \text{ as } x \\ \text{Inner terms} & N ::= X \mid \text{row_number}() \text{ over } (\text{order by } \vec{X}) \\ \text{Record terms} & R ::= \vec{N} \text{ as } \vec{\ell} \\ \text{Base terms} & X ::= x.\ell \mid c(\vec{X}) \end{array}$$

This fragment of SQL is almost isomorphic to the language we get if we apply let-insertion and record flattening to shredded queries. The only real difference is the use of `row_number` in place of `index`. Each instance of `index` in the body R of a subquery:

$$\text{for } (\overrightarrow{x \leftarrow t} \text{ where } X) \text{ return } R$$

is simulated by `row_number() over (order by $\overrightarrow{x.\ell}$)`, where $x_i : \langle \ell_{i,1} : A_{i,1}, \dots, \ell_{i,m_i} \rangle$, and:

$$\overrightarrow{x.\ell} = x_1.\ell_{1,1}, \dots, x_1.\ell_{1,m_1}, \dots, x_n.\ell_{n,1}, \dots, x_n.\ell_{n,m_n}$$

A possible concern is that `row_number` is non-deterministic. It computes row numbers ordered by the supplied columns, but if there is a tie, then it is free to order the equivalent rows in any order. However, such non-determinism is only observable if one computes a row number ordered by some columns of a table at the same time as inspecting other columns of that table. As we always order by *all* columns of all tables referenced from the current subquery, our use of `row_number` is always deterministic.

$$\begin{aligned}
\mathbf{L}(\biguplus_{i=1}^n C_i) &= \biguplus_{i=1}^n \mathbf{L}(C_i) \xrightarrow{\quad} \\
\mathbf{L}(C) &= \text{let } q = (\text{for } (\vec{G}_{\text{out}} \text{ where } X_{\text{out}}) \text{ return } \langle R_{\text{out}}, \text{index} \rangle) \text{ in for } (z \leftarrow q, \vec{G}_{\text{in}} \text{ where } \mathbf{L}_{\vec{y}}(X_{\text{in}})) \text{ return } \mathbf{L}_{\vec{y}}(N) \\
\text{where } \vec{G}_{\text{out}} &= \text{concat}(\text{init}(\text{gens } C)) \quad \vec{y} = \vec{t} = \vec{G}_{\text{out}} \quad \vec{G}_{\text{in}} = \text{last}(\text{gens } C) \quad N = \text{body } C \\
X_{\text{out}} &= \bigwedge \text{init}(\text{conds } C) \quad R_{\text{out}} = \langle \text{expand}(y_i, t_i) \rangle_{i=1}^n \quad X_{\text{in}} = \text{last}(\text{conds } C) \quad n = \text{length } \vec{G}_{\text{out}} \\
\mathbf{L}_{\vec{y}}(x.l) &= \begin{cases} x.l, & \text{if } x \notin \{y_1, \dots, y_n\} \\ z.1.i.l, & \text{if } x = y_i \end{cases} \quad \mathbf{L}_{\vec{y}}(\langle \ell_j = X_j \rangle_{j=1}^m) = \langle \ell_j = \mathbf{L}_{\vec{y}}(X_j) \rangle_{j=1}^m \quad \mathbf{L}(\text{out}) = z.2 \\
\mathbf{L}_{\vec{y}}(c(X_1, \dots, X_m)) &= c(\mathbf{L}_{\vec{y}}(X_1), \dots, \mathbf{L}_{\vec{y}}(X_m)) \quad \mathbf{L}_{\vec{y}}(a \diamond d) = \langle a, \mathbf{L}(d) \rangle \quad \mathbf{L}(\text{in}) = \text{index} \\
\text{expand}(x, t) &= \langle \ell_i = x.l_i \rangle_{i=1}^n \text{ where } \Sigma(t) = \text{Bag } \langle \ell : \vec{A} \rangle \quad \text{init } [x_i]_{i=1}^n = [x_i]_{i=1}^{n-1}, \dots, x_{n-1} \quad \text{last } [x_i]_{i=1}^n = x_n \\
\text{gens } (\text{for } (\vec{G} \text{ where } X) C) &= \vec{G} :: \text{gens } C \quad \text{conds } (\text{for } (\vec{G} \text{ where } X) C) = X :: \text{conds } C \quad \text{body } (\text{for } (\vec{G} \text{ where } X) C) = \text{body } C \\
\text{gens } (\text{return}^a N) &= [] \quad \text{conds } (\text{return}^a N) = [] \quad \text{body } (\text{return}^a N) = N
\end{aligned}$$

Figure 8. The let-insertion translation

$$\begin{aligned}
\mathcal{L}[\mathbf{L}] &= \mathcal{L}[\mathbf{L}]_\varepsilon & \mathcal{L}[\langle \ell_j = N_j \rangle_{j=1}^m]_{\rho, i} &= \langle \ell_j = \mathcal{L}[N_j]_{\rho, i} \rangle_{j=1}^m & \mathcal{L}[\langle \text{index} \rangle]_{\rho, i} &= i \\
\mathcal{L}[\biguplus_{j=1}^m C_j]_{\rho} &= \text{concat}(\mathcal{L}[C_j]_{\rho})_{j=1}^m & \mathcal{L}[\langle t \rangle]_{\rho} &= \langle t \rangle & \mathcal{L}[\langle q \rangle]_{\rho} &= \rho(q) \\
\mathcal{L}[\text{let } q = S_{\text{out}} \text{ in } S_{\text{in}}]_{\rho} &= \mathcal{L}[S_{\text{in}}]_{\rho[q \mapsto \mathcal{L}[S_{\text{out}}]_{\rho}]} \\
\mathcal{L}[\text{for } ([x_j \leftarrow u_j]_{j=1}^m \text{ where } X) \text{ return } N]_{\rho} &= \mathcal{L}[N]_{\rho[x_j \mapsto r_j]_{j=1}^m, i} \mid \langle i, \vec{r} \rangle \leftarrow \text{enum}([\vec{r} \mid [r_j \leftarrow \mathcal{L}[u_j]_{\rho}]_{j=1}^m, \mathcal{N}[X]_{\rho[x_j \mapsto r_j]_{j=1}^m}])
\end{aligned}$$

Figure 9. Semantics of let-inserted shredded queries

As an example, Q_2 becomes:

```

(with q as (select x.dept as i1_dept,
            row_number() over (order by x.dept) as i2
            from departments as x)
select a as i1_1, z.i2 as i1_2, y.name as i2_name, b as i2_tasks_1,
       row_number() over
       (order by z.i1_dept, z.i2, y.dept, y.employee, y.salary)
       as i2_tasks_2
from employees as y, q as z
where (z.i1_dept = y.dept ^
      (y.salary < 1000 v y.salary > 1000000)))
union all
(with q as (select x.dept as i1_dept,
            row_number() over (order by x.dept) as i2
            from departments as x)
select a as i1_1, z.i2 as i1_2, y.name as i2_name, d as i2_tasks_1,
       row_number() over
       (order by z.i1_dept, z.i2, y.dept, y.contact, y.client)
       as i2_tasks_2
from contacts as y, q as z
where (z.i1_dept = y.dept ^ y.client))

```

8. Variations and extensions

8.1 Empty, length and aggregation

Cooper's system includes the standard nested relational calculus primitive empty M for testing emptiness, which we can also add to our system. Given a term $M : \text{Bag } A$, empty M returns true iff M is empty. Normalisation proceeds on M as usual, with additional rules such as empty $\emptyset \rightsquigarrow \text{true}$. As the results themselves are not inspected, shredding need only generate the top level query for M . Ultimately, empty M is translated to a single SQL subquery. We can straightforwardly add primitives for length and for simple aggregation functions such as max in the same way.

8.2 Variant results

We have discussed variant types in the context of encoding shredded queries with natural indexes. We could easily support variant

types in their own right. However, if we disallow variants in results, then it is straightforward to normalise them away [28]. Indeed the current version of Links already does this [13]. Ulrich's Ferry-based extension to Links [24] also supports variants, including in the result type, taking the multiple query approach.

8.3 Higher-order results

Just as we relaxed the result type to admit nested results, we can relax it further to admit higher-order results. Again, this does not affect the correctness of normalisation. The resulting normal forms allow lambdas in the body of a comprehension, for instance:

$$\begin{aligned}
&\text{for } (x \leftarrow t) \text{ return } \lambda z.(x.l + z.l) \\
&\uplus \text{for } (x \leftarrow t', y \leftarrow t'') \text{ return } \lambda z.(x.l + y.l + z)
\end{aligned}$$

The goal is to encode lambdas in the results in such a way that a host programming language is able to reconstruct and apply them after running a query. This can be done straightforwardly using defunctionalisation [20], converting the lambdas into variants.

Ulrich [24] explored implementations of both closure conversion [1] and defunctionalisation in his work on integrating Ferry with Links, using them as an alternative to normalisation. He did not consider higher-order results or prove correctness.

8.4 List vs. set or bag semantics

Ultimately, we wish to integrate our shredding scheme with programming languages that manipulate lists. The normalisation scheme of Section 2 is based on a bag semantics. It is not sound for a list semantics because merging of conditionals into where clauses and hoisting of unions can lead to reordering of elements. The Links implementation addresses this issue by including additional ordering information in queries. We believe the same approach should apply to nested queries.

9. Related and future work

Understanding the relative expressiveness of different query languages, such as nested and flat relational query languages, is a classical topic in database theory, with many influential papers [18,

25, 28]. The Flat–Flat Theorem of Paredaens and Van Gucht [18] showed that nested relational queries are no more expressive than flat queries over flat relations. This result was strengthened by Wong [28], whose *Conservativity Theorem* showed that a query over inputs and outputs of nesting degree n does not require building intermediate data structures of greater degree. Wong’s proof gave a concrete algorithm for normalising nested relational queries over flat input data to a form that can be easily translated to SQL. Van den Bussche gave an alternative proof of the Flat–Flat Theorem by a simulation technique [25]. Recently, Cooper [7] extended Wong’s conservativity result to a higher-order nested relational calculus, in the context of the Links programming language [8]. This required a logical relations argument for collection types, using proof techniques that have only been developed recently; Cooper’s proof adapted the approach of Lindley and Stark [14].

These expressiveness results have influenced (and sometimes been influenced by) implementations. Wong’s work formed the basis for Kleisli [29], a commercially successful system used for biological data integration. Van den Bussche’s simulation has not been used directly as an implementation technique, but provides partial justification for Ferry’s loop-lifting algorithm. In addition, Ferry’s translation tries to handle lists, aggregation, grouping and ordering operations, as well as higher-order functions. Ferry’s approach has also recently been used as a backend for the Links system [24]. However, there is no proof of correctness for Ferry’s approach or its use within Links. Cooper’s approach to higher-order query normalisation has been incorporated into Links via a query keyword that identifies subexpressions expected to compile to an SQL query. Correct usage is enforced by an effect-based type system [7, 13].

Besides Cooper [7], several authors have recently considered higher-order query languages. Benedikt et al. [3] and Vu and Benedikt [26] study the complexity of containment, equivalence and evaluation for higher-order queries over flat relations. Higher-order features are also being added to XQuery 3.0.

Our work is also partly inspired by work on unnesting for nested data parallelism. Blleloch and Sabot [4] give a compilation scheme for NESL, a data-parallel language with nested lists; Suciu and Tannen [22] give an alternative scheme for a nested list calculus. This work may provide an alternative (and parallelisable) implementation strategy for Ferry’s list-based semantics [11].

The most obvious next step is to implement our algorithm and compare the performance of various indexing schemes (in terms of reliability, query generation and query evaluation) with Ferry’s loop-lifting technique and each other. In particular, we anticipate that it would be worthwhile to develop a schema-dependent *hybrid* scheme where natural indexes are used whenever possible.

10. Conclusion

Combining database programming with functional programming is challenging in part because of the limitations of flat database queries. Query shredding can help to bridge this gap. Although it is known from prior work that query shredding is possible in principle, and some implementations (such as Ferry) have tried to realise this, getting the details right is tricky. Our contribution is an alternative algorithm for shredding that handles higher-order queries over bags. More importantly, our approach gives a concrete and modular implementation strategy for shredding that is straightforward to prove correct, and should also be straightforward to extend and to incorporate into language-integrated query systems.

References

[1] A. W. Appel. *Compiling with continuations*. Cambridge University Press, New York, NY, USA, 1992.

[2] A. Beckmann. Exact bounds for lengths of reductions in typed λ -calculus. *Journal of Symbolic Logic*, 66, 2001.

[3] M. Benedikt, G. Puppis, and H. Vu. Positive higher-order queries. In *PODS*. ACM, 2010.

[4] G. E. Blleloch and G. W. Sabot. Compiling collection-oriented languages onto massively parallel computers. *J. Parallel Distrib. Comput.*, 8, February 1990.

[5] P. Buneman, L. Libkin, D. Suciu, V. Tannen, and L. Wong. Comprehension syntax. *SIGMOD Record*, 23, 1994.

[6] A. J. Chlipala. Ur: statically-typed metaprogramming with type-level record computation. In *PLDI*, 2010.

[7] E. Cooper. The script-writer’s dream: How to write great SQL in your own language, and be sure it will succeed. In *DBPL*, 2009.

[8] E. Cooper, S. Lindley, P. Wadler, and J. Yallop. Links: web programming without tiers. In *FMCO*, volume 4709 of *LNCS*, 2007.

[9] P. de Groote. On the strong normalisation of intuitionistic natural deduction with permutation-conversions. *Inf. Comput.*, 178(2), 2002.

[10] T. Grust, M. Mayr, J. Rittinger, and T. Schreiber. Ferry: Database-supported program execution. In *SIGMOD*, June 2009.

[11] T. Grust, J. Rittinger, and T. Schreiber. Avalanche-safe LINQ compilation. *PVLDB*, 3(1), 2010.

[12] S. Lindley. Extensional rewriting with sums. In *TLCA*, 2007.

[13] S. Lindley and J. Cheney. Row-based effect types for database integration. In *TLDI*, 2012.

[14] S. Lindley and I. Stark. Reducibility and \top -lifting for computation types. In *TLCA*, 2005.

[15] E. Meijer, B. Beckman, and G. M. Bierman. LINQ: reconciling object, relations and XML in the .NET framework. In *SIGMOD*, 2006.

[16] A. Ohori and K. Ueno. Making Standard ML a practical database programming language. In *ICFP*, 2011.

[17] C. Olston, B. Reed, U. Srivastava, R. Kumar, and A. Tomkins. Pig latin: a not-so-foreign language for data processing. In *SIGMOD*, 2008.

[18] J. Paredaens and D. Van Gucht. Converting nested algebra expressions into flat algebra expressions. *ACM Trans. Database Syst.*, 17(1), 1992.

[19] B. Pierce. *Types and Programming Languages*. MIT Press, 2002.

[20] J. C. Reynolds. Definitional interpreters for higher-order programming languages. *Higher-Order and Symbolic Computation*, 11(4), 1998.

[21] H. J. Schek and M. H. Scholl. The relational model with relation-valued attributes. *Inf. Syst.*, 11, April 1986.

[22] D. Suciu and V. Tannen. Efficient compilation of high-level data parallel algorithms. In *SPAA*. ACM, 1994.

[23] W. W. Tait. Intensional interpretations of functionals of finite type I. *Journal of Symbolic Logic*, 32(2), June 1967.

[24] A. Ulrich. A Ferry-based query backend for the Links programming language. Master’s thesis, University of Tübingen, 2011.

[25] J. Van den Bussche. Simulation of the nested relational algebra by the flat relational algebra, with an application to the complexity of evaluating powerset algebra expressions. *Theor. Comput. Sci.*, 254(1-2), 2001.

[26] H. Vu and M. Benedikt. Complexity of higher-order queries. In *ICDT*. ACM, 2011.

[27] P. Wadler. Comprehending monads. *Math. Struct. in Comp. Sci.*, 2(4), 1992.

[28] L. Wong. Normal forms and conservative extension properties for query languages over collection types. *J. Comput. Syst. Sci.*, 52(3), 1996.

[29] L. Wong. Kleisli, a functional query system. *J. Funct. Programming*, 10(1), 2000.

[30] XML Query and XSL Working Groups. XQuery 1.0: An XML query language. Available at <http://www.w3.org/TR/xquery/>, 2007.

A. Query normalisation

Following our previous work [13], we separate query normalisation into a rewriting phase and a type-directed structurally recursive function. Where we diverge from our previous work is that we extend the rewrite relation to hoist all conditionals up to the nearest conditional (in order to simplify the rest of the development), and the structurally recursive function is generalised to handle nested data. Thus normalisation can be divided into three stages.

- The first stage performs symbolic evaluation, that is, β -reduction and commuting conversions, flattening unnecessary nesting and eliminating higher-order functions.
- The second stage hoists all conditionals up to the nearest enclosing comprehension in order allow them to be converted to where clauses.
- The third and final stage hoists all unions up to the top-level, η -expands tables and variables, and turns all conditionals into where clauses.

Weak normalisation and strong normalisation Given a term M and a rewrite relation \rightsquigarrow , we write $M \not\rightsquigarrow$ if M is irreducible, that is no rewrite rules in \rightsquigarrow apply to M . The term M is said to be in *normal form*.

A term M is *weakly normalising* with respect to a rewrite relation \rightsquigarrow_r , or *r-WN*, if there exists a finite reduction sequence

$$M \rightsquigarrow M_1 \rightsquigarrow \dots \rightsquigarrow M_n \not\rightsquigarrow$$

A term M is *strongly normalising* with respect to a rewrite relation \rightsquigarrow_r , or *r-SN*, if every reduction sequence starting from M is finite. If M is *r-SN*, then we write $\max_r(M)$ for the maximum length of a reduction sequence starting from M .

A rewrite relation \rightsquigarrow_r is *weakly-normalising*, or *WN*, if all terms M are weakly-normalising with respect to \rightsquigarrow_r . Similarly, a rewrite relation \rightsquigarrow_r is *strongly-normalising*, or *SN*, if all terms M are strongly-normalising with respect to \rightsquigarrow_r .

We now describe each normalisation stage in turn.

A.1 Symbolic evaluation

The β -rules perform symbolic evaluation, including substituting argument values for function parameters, record field projection, conditionals where the test is known to be true or false, or iteration over singleton bags.

$$\begin{aligned} & (\lambda x. N) M \rightsquigarrow_c N[x := M] \\ & \langle \ell = M \rangle . \ell_i \rightsquigarrow_c M_i \\ & \text{if true then } M \text{ else } N \rightsquigarrow_c M \\ & \text{if false then } M \text{ else } N \rightsquigarrow_c N \\ & \text{for } (x \leftarrow \text{return } M) N \rightsquigarrow_c N[x := M] \end{aligned}$$

Fundamentally, β -rules always follow the same pattern. Each is associated with a particular type constructor T , and the left-hand side always consists of an *introduction* form for T inside an *elimination* form for T . For instance, in the case of functions (the first β -rule above), the introduction form is a lambda and the elimination form is an application. Applying a β -rule *eliminates* T in the sense that the introduction form from the left-hand side either no longer appears or has been replaced by a term of a simpler type on the right-hand side.

Each instance of a β -rule is associated with an *elimination frame*. An *elimination frame* is simply the elimination form with a designated *hole* $[\]$ that can be *plugged* with another expression. For instance, elimination frames for the function β -rule are of the form $E[\] = [\] M$. If we plug an introduction form $\lambda x. N$ into $E[\]$, written $E[\lambda x. N]$, then we obtain the left-hand side of the associated β -rule $(\lambda x. N) M$.

(This notion of an expression with a hole that can be filled by another expression is commonly used in rewriting and operational

semantics for programming languages [19, ch. 19]; here, it is not essential but helps cut down the number of explicit rules, and helps highlight commonality between rules.)

The elimination frames of λ_{NRC} are as follows.

$$E[\] ::= [\] M \mid [\].\ell \mid \text{if } [\] \text{ then } M \text{ else } N \mid \text{for } (x \leftarrow [\]) N$$

The following rules express that comprehensions, conditionals, empty bag constructors, and unions can always be *hoisted* out of the above elimination frames. In the literature such rules are often called *commuting conversions*. They are necessary in order to expose all possible β -reductions. For instance,

$$(\text{if } M \text{ then } \langle \ell = N \rangle \text{ else } N') . \ell$$

cannot β -reduce, but if $M \text{ then } \langle \ell = N \rangle . \ell \text{ else } N' . \ell$ can.

$$\begin{aligned} E[\text{for } (x \leftarrow M) N] & \rightsquigarrow_c \text{for } (x \leftarrow M) E[N] \\ E[\text{if } L \text{ then } M \text{ else } N] & \rightsquigarrow_c \text{if } L \text{ then } E[M] \text{ else } E[N] \\ E[\emptyset] & \rightsquigarrow_c \emptyset \\ E[M_1 \uplus M_2] & \rightsquigarrow_c E[M_1] \uplus E[M_2] \end{aligned}$$

For example:

$$(\text{if } L \text{ then } M_1 \text{ else } M_2) M \rightsquigarrow_c \text{if } L \text{ then } M_1 M \text{ else } M_2 M$$

Note that some combinations of elimination frames and rewrite rule are impossible in a well-typed term, such as if \emptyset then M else N .

Next we prove that \rightsquigarrow_c is strongly normalising. The proof is based on our previous proof of strong normalisation for simply-typed λ -calculus with sums [12], which generalises the $\top\top$ -lifting approach [14], which in turn extends Tait's proof of strong normalisation for simply-typed λ -calculus [23].

Frame stacks

$$\begin{aligned} (\text{frame stacks}) \quad S & ::= Id \mid S \circ E \\ (\text{stack length}) \quad |Id| & = 0 \\ & |S \circ E| = |S| + 1 \\ (\text{plugging}) \quad Id[M] & = M \\ (S \circ E)[M] & = S[(E[M])] \end{aligned}$$

Following our previous work [12] we assume variables are annotated with types. We write $A \multimap B$ for the type of frame stack S , if $S[M] : B$ for all terms $M : A$.

Frame stack reduction

$$S \rightsquigarrow_c S' \stackrel{\text{def}}{\iff} \forall M. S[M] \rightsquigarrow_c S'[M] \iff S[x] \rightsquigarrow_c S'[x]$$

Frame stacks are closed under reduction. A frame stack S is *c-strongly normalising*, or *c-SN*, if all reduction sequences starting from S are finite.

Lemma 12. *If $S \rightsquigarrow_c S'$, for frame stacks S, S' , then $|S'| \leq |S|$.*

Proof. Induction on the structure of S . □

Reducibility We define *reducibility* as follows:

- Id is reducible.
- $S \circ ([\] N) : (A \rightarrow B) \multimap C$ is reducible if S and N are reducible.
- $S \circ ([\].\ell) : (A \times B) \multimap C$ is reducible if S is reducible.
- $S : \text{Bag } A \multimap C$ is reducible if $S[\text{return } M]$ is *c-SN* for all reducible $M : A$.
- $S : \text{Bool} \multimap C$ is reducible if $S[\text{true}]$ is *c-SN* and $S[\text{false}]$ is *c-SN*.
- $M : A$ is reducible if $S[M]$ is *c-SN* for all reducible $S : A \multimap C$.

Lemma 13. *If $M : A$ is reducible then M is *c-SN*.*

Proof. Follows immediately from reducibility of Id and the definition of reducibility on terms. \square

Lemma 14. $x : A$ is reducible.

Proof. By induction on A using Lemma 12 and Lemma 13. \square

Corollary 15. If $S : A \multimap C$ is reducible then S is c -SN.

Each type constructor has an associated β -rule. Each β -rule gives rise to an SN-closure property.

Lemma 16 (SN-closure).

(\rightarrow) If $S[M[x := N]]$ and N are c -SN then $S[(\lambda x.M) N]$ is c -SN.

($\langle \rangle$) If \vec{M} , are c -SN then $S[\langle \ell = \vec{M} \rangle . \ell_i]$ is c -SN.

(Bag \rightarrow) If $S[N[x := M]]$ and M are c -SN then $S[\text{for } (x \leftarrow \text{return } M) N]$ is c -SN.

(Bool) If $S[N]$ and $S[N']$ are c -SN then $S[\text{if true then } N \text{ else } N']$ is c -SN and $S[\text{if false then } N \text{ else } N']$ is c -SN.

Proof.

(\rightarrow): By induction on $\max_c(S) + \max_c(M) + \max_c(N)$.

($\langle \rangle$): By induction on $\max_c S + (\sum_{i=1}^n \max_c(M_i)) + \max_c(N) + (\sum_{i=1}^n \max_c(M'_i))$.

(Bag \rightarrow): By induction on $|S| + \max_c(S[N[x := M]]) + \max_c(M)$.

(Bool): By induction on $|S| + \max_c(S[N]) + \max_c(S[N'])$. \square

Now we obtain reducibility-closure properties for each type constructor.

Lemma 17 (reducibility-closure).

(\rightarrow) If $M[x := N]$ is reducible for all reducible N , then $\lambda x.M$ is reducible.

($\langle \rangle$) If \vec{M} are reducible, then $\langle \ell = \vec{M} \rangle$ is reducible.

(Bag) If M is reducible, $N[x := M']$ is reducible for all reducible M' , then for $(x \leftarrow M) N$ is reducible.

(Bool) If M, N, N' are reducible then if M then N else, N' is reducible.

Proof. Each property follows from the corresponding part of Lemma 16 using Lemma 13 and Corollary 15. \square

We also require additional closure properties for the empty bag and union constructs.

Lemma 18 (reducibility-closure II).

(\emptyset) The empty bag \emptyset is reducible.

(\uplus) If M, N are reducible, then $M \uplus N$ is reducible.

Proof.

(\emptyset): Suppose $S : \text{Bag } A \multimap C$ is reducible. We need to prove that $S[\emptyset]$ is c -SN. The proof is by induction on $|S| + \max_c(S)$. The only interesting case is hoisting the empty bag out of a bag elimination frame, which simply decreases the size of the frame stack by 1.

(\uplus): Suppose $M, N : \text{Bag } A$, and $S : \text{Bag } A \multimap C$ are reducible. We need to show that $S[M \uplus N]$ is c -SN. The proof is by induction on $|S| + \max_c(S[M]) + \max_c(S[N])$. The only interesting case is hoisting the union out of a bag elimination frame, which again decreases the size of the frame stack by 1, whilst leaving the other components of the induction measure unchanged. \square

Theorem 19. Let M be any term. Suppose $x_1 : A_1, \dots, x_n : A_n$ includes all the free variables of M . If $N_1 : A_1, \dots, N_n : A_n$ are reducible then $M[\vec{x} := \vec{N}]$ is reducible.

Proof. By induction on the structure of terms using Lemma 17 and Lemma 18. \square

Theorem 20 (strong normalisation). The relation \rightsquigarrow_c is strongly normalising.

Proof. Let M be a term with free variables \vec{x} . By Lemma 14, \vec{x} are reducible. Hence, by Theorem 19, M is c -SN. \square

It is well known that β -reduction in simply-typed λ -calculus has non-elementary complexity in the worst case [2]. The relation \rightsquigarrow_c includes β -reduction, so it must be at least as bad (we conjecture that it has the same asymptotic complexity, as \rightsquigarrow_c can be reduced to β -reduction on simply-typed λ -calculus via a CPS translation following de Groote [9]). However, we believe that the asymptotic complexity is unlikely to pose a problem in practice, as the kind of higher-order code that exhibits worst-case behaviour is rare. It has been our experience with Links that query normalisation time is almost always dominated by SQL execution time.

A.2 If hoisting

To hoist conditionals (if-expressions) out of constant applications, records, unions, and singleton bag constructors, we define if-hoisting frames as follows:

$$F[\] ::= c(\vec{M}, [\], \vec{N}) \mid \langle \ell' = \vec{M}, \ell = [\], \ell'' = \vec{N} \rangle \\ \mid [\] \uplus N \mid M \uplus [\] \mid \text{return } [\]$$

The if-hoisting rule says that if an expression contains an if-hoisting frame around a conditional, then we can lift the conditional up and push the frame into both branches:

$$F[\text{if } L \text{ then } M \text{ else } N] \rightsquigarrow_h \text{if } L \text{ then } F[M] \text{ else } F[N]$$

We write $\text{size}(M)$ for the size of M , as in the total number of syntax constructors in M .

Lemma 21.

1. If M, N, N' are h -SN then if M then N else N' is h -SN.
2. If M, N are h -SN then $M \uplus N$ is h -SN.
3. If \vec{M} are h -SN then $c(\vec{M})$ is h -SN.
4. If \vec{M} are h -SN then $\langle \ell = \vec{M} \rangle$ is h -SN.

Proof.

1: By induction on $\langle \max_h(M), \text{size}(M), \max_h(N) + \max_h(N'), \text{size}(N) + \text{size}(N') \rangle$.

2: By induction on $\langle \max_h(M) + \max_h(N), \text{size}_h(M) + \text{size}_h(N) \rangle$ using (1).

3 and 4: By induction on $\langle \sum_{i=1}^n \max_h(M_i), \sum_{i=1}^n \text{size}(M_i) \rangle$ using (1). \square

Proposition 22. The relation \rightsquigarrow_h is strongly normalising.

Proof. By induction on the structure of terms using Lemma 21. \square

A.3 The query normalisation function

The following definition of the function $norm$ generalises the normalisation algorithm from our previous work [13].

$$norm_A(M) = \langle N \rangle_A$$

where $M \rightsquigarrow_c^* M' \rightsquigarrow_h^* N$ with $M' \not\rightsquigarrow_c, N \not\rightsquigarrow_h$ and:

$$\begin{aligned} \langle M \rangle_O &= M \\ \langle M \rangle_{\langle \ell_i : A_i \rangle_{i=1}^n} &= \langle \ell_i = \mathcal{F}(M)_{A_i, \ell_i} \rangle_{i=1}^n \\ \langle M \rangle_{\text{Bag } A} &= \biguplus (\mathcal{B}(M)_{A, [\cdot], \text{true}}^*) \\ \mathcal{B}(\text{return } M)_{A, \vec{G}, L}^* &= [\text{for } (\vec{G} \text{ where } L) \text{ return } \langle M \rangle_A] \\ \mathcal{B}(\text{for } (x \leftarrow t) M)_{A, \vec{G}, L}^* &= \mathcal{B}(M)_{A, \vec{G} ++ [x \leftarrow t], L}^* \\ \mathcal{B}(\text{table } t)_{A, \vec{G}, L}^* &= \mathcal{B}(x)_{A, \vec{G} ++ [x \leftarrow t], L}^* \quad (x \text{ fresh}) \\ \mathcal{B}(\emptyset)_{A, \vec{G}, L}^* &= [] \\ \mathcal{B}(M \uplus N)_{A, \vec{G}, L}^* &= \mathcal{B}(M)_{A, \vec{G}, L}^* ++ \mathcal{B}(N)_{A, \vec{G}, L}^* \\ \mathcal{B}(\text{if } L' M N)_{A, \vec{G}, L}^* &= \mathcal{B}(M)_{A, \vec{G}, L \wedge L'}^* \\ &\quad ++ \mathcal{B}(N)_{A, \vec{G}, L \wedge \neg L'}^* \\ \mathcal{F}(x)_{A, \ell_i} &= \langle x, \ell_i \rangle_A \\ \mathcal{F}(\langle \ell_i = M_i \rangle_{i=1}^n)_{A, \ell_i} &= \langle M_i \rangle_A \end{aligned}$$

Strictly speaking, in order for the above definition of $norm_A$ to make sense we need the two rewrite relations to be confluent. It is easily verified that the relation \rightsquigarrow_c is locally confluent, and hence by strong normalisation and Newman's Lemma it is confluent. The relation \rightsquigarrow_h is not confluent as the ordering of hoisting determines the final order in which booleans are eliminated. However, it is easily seen to be confluent modulo reordering of conditionals, which in turn means that $norm_A$ is well-defined if we identify terms modulo commutativity of conjunction, which is perfectly reasonable given that conjunction is indeed commutative.

Theorem 23. *The function $norm_A$ terminates.*

Proof. The result follows immediately from strong normalisation of \rightsquigarrow_c and \rightsquigarrow_h , and the fact that the functions $\langle - \rangle$, $\mathcal{B}(-)^*$, and $\mathcal{F}(-)$ are structurally recursive (modulo expanding out the right-hand-side of the definition of $\mathcal{B}(\text{table } t)_{A, \vec{G}, L}^*$). \square

B. Typing rules for shredded terms and values

The typing rules for shredded terms and values are shown in Figures 10, 11 and 12.

C. Proof of correctness of shredding

We begin with a simple lemma that states that the inner shredding function $\llbracket - \rrbracket$ commutes with the semantics.

Lemma 24. $S[\llbracket M \rrbracket_a]_{\rho, \iota} = \llbracket \mathcal{A}[M]_{\rho, \iota} \rrbracket_{a \diamond \iota}$

Proof. By induction on the structure of M .

Case X :

$$\begin{aligned} &S[\llbracket X \rrbracket_a]_{\rho, \iota} \\ &= S[X]_{\rho, \iota} \\ &= \mathcal{A}[X]_{\rho, \iota} \\ &= \llbracket \mathcal{A}[X]_{\rho, \iota} \rrbracket_{a \diamond \iota} \end{aligned}$$

Case $\langle \ell = M \rangle$:

$$\begin{aligned} &S[\llbracket \langle \ell_i = M_i \rangle_{i=1}^n \rrbracket_a]_{\rho, \iota} \\ &= S[\langle \ell_i = \llbracket M_i \rrbracket_a \rangle_{i=1}^n]_{\rho, \iota} \\ &= \langle \ell_i = S[\llbracket M_i \rrbracket_a]_{\rho, \iota} \rangle_{i=1}^n \\ &\quad (\text{Induction hypothesis}) \\ &= \langle \ell_i = \llbracket \mathcal{A}[M_i]_{\rho, \iota} \rrbracket_{a \diamond \iota} \rangle_{i=1}^n \\ &= \llbracket \langle \ell_i = \mathcal{A}[M_i]_{\rho, \iota} \rangle_{i=1}^n \rrbracket_{a \diamond \iota} \\ &= \llbracket \mathcal{A}[\langle \ell_i = M_i \rangle_{i=1}^n]_{\rho, \iota} \rrbracket_{a \diamond \iota} \end{aligned}$$

Case L :

$$\begin{aligned} &S[\llbracket L \rrbracket_a]_{\rho, \iota} \\ &= S[a \diamond \text{in}]_{\rho, \iota} \\ &= a \diamond \iota \\ &= \llbracket \mathcal{A}[L]_{\rho, \iota} \rrbracket_{a \diamond \iota} \end{aligned}$$

\square

The first part of the following lemma allows us to run a shredded query by concatenating the results of running the shreds of its component comprehensions. Similarly, the second part allows us to shred the results of running a nested query by concatenating the shreds of the results of running its component comprehensions.

Lemma 25.

1. $S[\llbracket \biguplus_{i=1}^n C_i \rrbracket_{a, p}^*]_{\rho, \iota} = \text{concat}(\{S[\llbracket C_i \rrbracket_{a, p}^*]_{\rho, \iota}\}_{i=1}^n)$
2. $\llbracket \mathcal{A}[\biguplus_{i=1}^n C_i]_{\rho, \iota} \rrbracket_{a \diamond \iota, p}^* = \text{concat}(\{\llbracket \mathcal{A}[C_i]_{\rho, \iota} \rrbracket_{a \diamond \iota, p}^* \rrbracket_{i=1}^n\})$

Proof.

1.

$$\begin{aligned} &S[\llbracket \biguplus_{i=1}^n C_i \rrbracket_{a, p}^*]_{\rho, \iota} \\ &= S[\llbracket \biguplus \text{concat}(\{\llbracket C_i \rrbracket_{a, p}^* \rrbracket_{i=1}^n\}) \rrbracket_{\rho, \iota}]_{\rho, \iota} \\ &= S[\llbracket \biguplus \text{concat}(\{\llbracket C \rrbracket_{a, p}^* \mid C \leftarrow \vec{C}\}) \rrbracket_{\rho, \iota}]_{\rho, \iota} \\ &= \text{concat}(\text{concat}(\{S[\llbracket C' \rrbracket_{\rho, \iota} \mid C \leftarrow \vec{C}, C' \leftarrow \llbracket C \rrbracket_{a, p}^* \rrbracket_{\rho, \iota}\})\}) \\ &= \text{concat}(\{S[\llbracket \biguplus \llbracket C \rrbracket_{a, p}^* \rrbracket_{\rho, \iota} \mid C \leftarrow \vec{C}\}]\}) \\ &= \text{concat}(\{S[\llbracket \biguplus \llbracket C_i \rrbracket_{a, p}^* \rrbracket_{\rho, \iota} \rrbracket_{i=1}^n]\}) \end{aligned}$$
2.

$$\begin{aligned} &\llbracket \mathcal{A}[\biguplus_{i=1}^n C_i]_{\rho, \iota} \rrbracket_{a \diamond \iota, p}^* \\ &= \llbracket \text{concat}(\{\mathcal{A}[C_i]_{\rho, \iota}\}_{i=1}^n) \rrbracket_{a \diamond \iota, p}^* \\ &= \llbracket \text{concat}(\{\mathcal{A}[C]_{\rho, \iota} \mid C \leftarrow \vec{C}\}) \rrbracket_{a \diamond \iota, p}^* \\ &= \text{concat}(\text{concat}(\{\llbracket s \rrbracket_{a \diamond \iota, p}^* \mid C \leftarrow \vec{C}, s \leftarrow \mathcal{A}[C]_{\rho, \iota}\})\}) \\ &= \text{concat}(\{\llbracket \mathcal{A}[C]_{\rho, \iota} \rrbracket_{a \diamond \iota, p}^* \mid C \leftarrow \vec{C}\}) \\ &= \text{concat}(\{\llbracket \mathcal{A}[C_i]_{\rho, \iota} \rrbracket_{a \diamond \iota, p}^* \rrbracket_{i=1}^n\}) \end{aligned}$$

$$\begin{array}{c}
\text{VAR} \\
\frac{}{\Gamma, x : \langle \ell : F \rangle \vdash x : \langle \ell : F \rangle} \\
\\
\text{PROJECT} \\
\frac{\Gamma \vdash x : \langle \ell_i : A_i \rangle_{i=1}^n}{\Gamma \vdash x.\ell_j : A_j} \\
\\
\text{INDEX} \\
\frac{}{\Gamma \vdash a \diamond d : Index} \\
\\
\text{SINGLETON} \\
\frac{\Gamma \vdash I : Index \quad \Gamma \vdash N : F}{\Gamma \vdash \text{return}^a \langle I, N \rangle : \text{Bag} \langle Index, F \rangle} \\
\\
\text{UNION} \\
\frac{\Gamma \vdash C_i : \text{Bag} \langle Index, F \rangle_{i=1}^n}{\Gamma \vdash \biguplus \vec{C} : \text{Bag} \langle Index, F \rangle} \\
\\
\text{CONSTANT} \\
\frac{\Sigma(c) = \langle O_1, \dots, O_n \rangle \rightarrow O' \quad [\Gamma \vdash X_i : O_i]_{i=1}^n}{\Gamma \vdash c \langle X_1, \dots, X_n \rangle : O'} \\
\\
\text{RECORD} \\
\frac{[\Gamma \vdash N_i : A_i]_{i=1}^n}{\Gamma \vdash \langle \ell_i = N_i \rangle_{i=1}^n : \langle \ell_i = A_i \rangle_{i=1}^n} \\
\\
\text{FOR} \\
\frac{[\Sigma(t_i) = \text{Bag} A_i]_{i=1}^n \quad \Gamma, [x_i : A_i]_{i=1}^n \vdash X : Bool \quad \Gamma, [x_i : A_i]_{i=1}^n \vdash C : \text{Bag} \langle Index, F \rangle}{\Gamma \vdash \text{for}([x_i \leftarrow t_i]_{i=1}^n \text{ where } X) C : \text{Bag} \langle Index, F \rangle}
\end{array}$$

Figure 10. Typing rules for shredded terms

$$\begin{array}{c}
\text{RESULT} \\
\frac{[\vdash v_i : A_i]_{i=1}^n \quad [\vdash J : Index]_{i=1}^n}{\vdash [v_i @ J_i]_{i=1}^n : \text{Bag } A} \\
\\
\text{RECORD} \\
\frac{[\vdash v_i : A_i]_{i=1}^n}{\vdash \langle \ell_i = v_i \rangle_{i=1}^n : \langle \ell_i = A_i \rangle_{i=1}^n} \\
\\
\text{CONSTANT} \\
\frac{\Sigma(c) = O}{\vdash c : O}
\end{array}$$

Figure 11. Typing rules for indexed nested values

$$\begin{array}{c}
\text{RESULT} \\
\frac{[\vdash I_i : Index]_{i=1}^n \quad [\vdash w_i : F]_{i=1}^n \quad [\vdash J_i : Index]_{i=1}^n}{\vdash \langle I_1, w_1 @ J_1 \rangle, \dots, \langle I_n, w_n @ J_n \rangle : \text{Bag} \langle Index, F \rangle} \\
\\
\text{CONSTANT} \\
\frac{\Sigma(c) = O}{\vdash c : O} \\
\\
\text{RECORD} \\
\frac{[\vdash n_i : F_i]_{i=1}^n}{\vdash \langle \ell_i = w_i \rangle_{i=1}^n : \langle \ell_i = F_i \rangle_{i=1}^n}
\end{array}$$

Figure 12. Typing rules for shredded values

□ **Case $\downarrow.\epsilon$:**

We are now in a position to prove that the outer shredding function $\llbracket - \rrbracket$ commutes with the semantics.

Lemma 26. $S[\llbracket \llbracket L \rrbracket_p \rrbracket_{\rho,1}] = \llbracket \mathcal{A}[\llbracket L \rrbracket_{\rho,1}] \rrbracket_p$

Proof. We prove the following:

1. $S[\llbracket \llbracket L \rrbracket_p \rrbracket_{\rho,1}] = \llbracket \mathcal{A}[\llbracket L \rrbracket_{\rho,1}] \rrbracket_p$
 2. $S[\llbracket \llbracket C \rrbracket_{a,p}^* \rrbracket_{\rho,t}] = \llbracket \mathcal{A}[\llbracket C \rrbracket_{\rho,t}] \rrbracket_{a \diamond \iota, p}^*$
 3. $S[\llbracket \llbracket M \rrbracket_{a,p}^* \rrbracket_{\rho,t}] = \llbracket \mathcal{A}[\llbracket M \rrbracket_{\rho,t}] \rrbracket_{a \diamond \iota, p}^*$
- The first equation is the result we require. Observe that it follows from (2) and Lemma 25. We now proceed to prove equations (2) and (3) by mutual induction on the structure of p .

There are only two cases for (2), as the $\ell_i.p$ case cannot apply.

Case ϵ :

$$\begin{aligned}
& S[\llbracket \llbracket \text{for}(\overrightarrow{x \leftarrow t} \text{ where } X) \text{return}^b M \rrbracket_{a,\epsilon}^* \rrbracket_{\rho,\iota} \\
&= \text{(Definition of } S[\llbracket - \rrbracket]) \\
& S[\llbracket \text{for}(\overrightarrow{x \leftarrow t} \text{ where } X) \text{return}^b \langle a \diamond \iota, \llbracket M \rrbracket_b \rangle \rrbracket_{\rho,\iota} \\
&= \text{(Definition of } S[\llbracket - \rrbracket]) \\
& \llbracket \langle a \diamond \iota, S[\llbracket \llbracket M \rrbracket_b \rrbracket_{\rho[\overrightarrow{x \rightarrow t}], \iota, i}] \rrbracket_{a \diamond \iota, i} \rrbracket_{a \diamond \iota, \iota} \\
& \quad | \langle i, \vec{v} \rangle \leftarrow \text{enum}([\vec{v} \mid v \leftarrow \llbracket t \rrbracket], S[\llbracket X \rrbracket_{\rho[\overrightarrow{x \rightarrow t}]}]) \\
&= \text{(Lemma 24)} \\
& \llbracket \langle a \diamond \iota, \llbracket \mathcal{A}[\llbracket M \rrbracket_{\rho[\overrightarrow{x \rightarrow t}], \iota, i}] \rrbracket_{b \diamond \iota, i} \rrbracket_{b \diamond \iota, \iota} \rrbracket_{a \diamond \iota, \iota} \\
& \quad | \langle i, \vec{v} \rangle \leftarrow \text{enum}([\vec{v} \mid v \leftarrow \llbracket t \rrbracket], S[\llbracket X \rrbracket_{\rho[\overrightarrow{x \rightarrow t}]}]) \\
&= \text{(Definition of } \llbracket - \rrbracket^*) \\
& \llbracket \llbracket \mathcal{A}[\llbracket M \rrbracket_{\rho[\overrightarrow{x \rightarrow t}], \iota, i}] \rrbracket_{b \diamond \iota, i} \rrbracket_{a \diamond \iota, \iota} \\
& \quad | \langle i, \vec{v} \rangle \leftarrow \text{enum}([\vec{v} \mid v \leftarrow \llbracket t \rrbracket], S[\llbracket X \rrbracket_{\rho[\overrightarrow{x \rightarrow t}]}]) \rrbracket_{a \diamond \iota, \epsilon}^* \\
&= \text{(Definition of } \mathcal{A}[\llbracket - \rrbracket]) \\
& \llbracket \mathcal{A}[\llbracket \text{for}(\overrightarrow{x \leftarrow t} \text{ where } X) \text{return}^b M \rrbracket_{\rho,\iota}] \rrbracket_{a \diamond \iota, \epsilon}^*
\end{aligned}$$

$$\begin{aligned}
& S[\llbracket \llbracket \text{for}(\overrightarrow{x \leftarrow t} \text{ where } X) \text{return}^b M \rrbracket_{a,\downarrow.p}^* \rrbracket_{\rho,\iota} \\
&= \text{(Definition of } \llbracket - \rrbracket^*) \\
& S[\llbracket \llbracket \text{for}(\overrightarrow{x \leftarrow t} \text{ where } X) C \mid C \leftarrow \llbracket M \rrbracket_{b,p}^* \rrbracket_{\rho,\iota} \\
&= \text{(Definition of } S[\llbracket - \rrbracket]) \\
& \text{concat}([\text{concat} \\
& \quad ([S[\llbracket C \rrbracket_{\rho[\overrightarrow{x \rightarrow t}], \iota, i}] \\
& \quad \quad | \langle i, \vec{v} \rangle \leftarrow \text{enum}([\vec{v} \mid v \leftarrow \llbracket t \rrbracket], S[\llbracket X \rrbracket_{\rho[\overrightarrow{x \rightarrow t}]}]) \\
& \quad \quad | C \leftarrow \llbracket M \rrbracket_{b,p}^*]) \\
&= \text{(Definition of } S[\llbracket - \rrbracket]) \\
& \text{concat} \\
& \quad ([S[\llbracket \llbracket M \rrbracket_{b,p}^* \rrbracket_{\rho[\overrightarrow{x \rightarrow t}], \iota, i}] \\
& \quad \quad | \langle i, \vec{v} \rangle \leftarrow \text{enum}([\vec{v} \mid v \leftarrow \llbracket t \rrbracket], S[\llbracket X \rrbracket_{\rho[\overrightarrow{x \rightarrow t}]}]) \\
&= \text{(Induction hypothesis (3))} \\
& \text{concat} \\
& \quad ([\llbracket \mathcal{A}[\llbracket M \rrbracket_{\rho[\overrightarrow{x \rightarrow t}], \iota, i}] \rrbracket_{b,p}^* \\
& \quad \quad | \langle i, \vec{v} \rangle \leftarrow \text{enum}([\vec{v} \mid v \leftarrow \llbracket t \rrbracket], S[\llbracket X \rrbracket_{\rho[\overrightarrow{x \rightarrow t}]}]) \\
&= \text{(Definitions of } \mathcal{A}[\llbracket - \rrbracket] \text{ and } \llbracket - \rrbracket^*) \\
& \llbracket \mathcal{A}[\llbracket \text{for}(\overrightarrow{x \leftarrow t} \text{ where } X) \text{return}^b M \rrbracket_{\rho,\iota}] \rrbracket_{a \diamond \iota, \downarrow.p}^*
\end{aligned}$$

There are three cases for (3).

Cases ϵ and $\downarrow.\epsilon$: follow from (2) by applying the two parts of Lemma 25 to the left and right-hand side respectively.

Case $\ell_i.\epsilon$:

$$\begin{aligned}
& \mathcal{S}[\llbracket \langle \ell = \overrightarrow{M} \rangle \rrbracket_{a,\ell_i,p}^*]_{\rho,\iota} \\
&= \text{(Definition of } \llbracket - \rrbracket^* \text{)} \\
&= \mathcal{S}[\llbracket M_i \rrbracket_{a,p}^*]_{\rho,\iota} \\
&= \text{(Induction hypothesis (3))} \\
&= \llbracket \mathcal{A}[M_i]_{\rho,\iota} \rrbracket_{a \diamond \iota, p}^* \\
&= \text{(Definition of } \mathcal{A}[\llbracket - \rrbracket] \text{)} \\
&= \llbracket \mathcal{A}[\langle \ell = \overrightarrow{M} \rangle]_{\rho,\iota} \rrbracket_{a \diamond \iota, \ell_i, p}^*
\end{aligned}$$

□

We now lift Lemma 26 to shredded packages.

Proof of Theorem 5. We need to show that if $\vdash L : \text{Bag } A$ then

$$\mathcal{H}[\llbracket L \rrbracket_A] = \text{shred}_{\mathcal{A}[\llbracket L \rrbracket]}(\text{Bag } A)$$

This is straightforward by induction on A , using Lemma 26. □

We have proved that shredding commutes with the semantics. It remains to show that stitching after shredding is the identity on index-annotated nested results. We need two auxiliary notions: the *descendant* of a value at a path, and the *indexes* at the end of a path (which must be unique in order for stitching to work).

We define $\llbracket v \rrbracket_{p,w}$, the *descendant* of a value v at path p with respect to inner value w as follows.

$$\begin{aligned}
\llbracket v \rrbracket_{p,w} &= \llbracket v \rrbracket_{\top \diamond 1, p, w} \\
\llbracket v \rrbracket_{J,p,c} &= c \\
\llbracket v \rrbracket_{J,p, \langle \ell_i = w_i \rangle_{i=1}^n} &= \langle \ell_i = \llbracket v \rrbracket_{J,p, \ell_i, w_i} \rangle_{i=1}^n \\
\llbracket s \rrbracket_{J, \epsilon, I} &= \begin{cases} s, & \text{if } J = I \\ [], & \text{if } J \neq I \end{cases} \\
\llbracket [v_i @ J_i]_{i=1}^n \rrbracket_{J, \downarrow, p, I} &= \text{concat}(\llbracket v_i \rrbracket_{J_i, p, I} \rangle_{i=1}^n) \\
\llbracket \langle \ell_i = v_i \rangle_{i=1}^n \rrbracket_{J, \ell_i, p, I} &= \llbracket v_i \rrbracket_{J, p, I}
\end{aligned}$$

Essentially, this extracts the part of v that corresponds to w . The inner value w allows us to specify a particular descendant as an index, or nested record of indexes; for uniformity it may also contain constants.

We can now formulate the following crucial technical lemma, which states that given the descendants of a result v at path $p.\downarrow$ and the shredded values of v at path p we can stitch them together to form the descendants at path p .

Lemma 27 (key lemma). *If v is well-indexed and $\vdash \llbracket v \rrbracket_{J,p,I} : \text{Bag } A$, then*

$$\llbracket v \rrbracket_{J,p,I} = \llbracket \llbracket v \rrbracket_{J,p,\downarrow,w} @ I_{in} \mid \langle I_{out}, w \rangle @ I_{in} \leftarrow \llbracket v \rrbracket_{J,p}^*, I_{out} = I \rrbracket$$

Proof. First we strengthen the induction hypothesis to account for records. The generalised induction hypothesis is as follows.

If v is well-indexed and $\vdash \llbracket v \rrbracket_{J',p,\vec{\ell},I} : \text{Bag } A$, then

$$\begin{aligned}
& \llbracket \llbracket v \rrbracket_{J',p,\downarrow,\vec{\ell},w} @ I_{in} \mid \langle I_{out}, w \rangle @ I_{in} \leftarrow \llbracket v \rrbracket_{J',p,\vec{\ell}}^*, I_{out} = I \rrbracket \\
&= \llbracket v' . \vec{\ell} @ I_{in} \mid v' @ I_{in} \leftarrow \llbracket v \rrbracket_{J',p,I} \rrbracket
\end{aligned}$$

The proof now proceeds by induction on the structure of A and side-induction on the structure of p .

Case O :

Subcase ϵ : If $J' \neq I$ then both sides are empty lists. Suppose that $J' = I$.

$$\begin{aligned}
& \llbracket \llbracket s \rrbracket_{J',\downarrow,\vec{\ell},c} @ I_{in} \mid \langle I_{out}, c \rangle @ I_{in} \leftarrow \llbracket s \rrbracket_{J',\vec{\ell}}^*, I_{out} = I \rrbracket \\
&= \text{(Definition of } \llbracket - \rrbracket^* \text{)} \\
&= \llbracket c @ I_{in} \mid \langle I_{out}, c \rangle @ I_{in} \leftarrow \llbracket s \rrbracket_{J',\vec{\ell}}^*, I_{out} = I \rrbracket \\
&= (J' = I) \\
&= s \\
&= \text{(Definition of } \llbracket - \rrbracket^* \text{)} \\
&= \llbracket v' . \vec{\ell} @ I_{in} \mid v' @ I_{in} \leftarrow \llbracket s \rrbracket_{J',\epsilon,I} \rrbracket
\end{aligned}$$

Subcase $\ell_i.p$:

$$\begin{aligned}
& \llbracket \llbracket \langle \ell = \overrightarrow{v} \rangle \rrbracket_{J',\ell_i.p,\downarrow,\vec{\ell},c} @ I_{in} \mid \langle I_{out}, c @ I_{in} \rangle @ I_{in} \leftarrow \llbracket \langle \ell = \overrightarrow{v} \rangle \rrbracket_{J',\ell_i.p,\vec{\ell}}^*, I_{out} = I \rrbracket \\
&= \text{(Definition of } \llbracket - \rrbracket^* \text{)} \\
&= \llbracket c @ I_{in} \mid \langle I_{out}, c @ I_{in} \rangle @ I_{in} \leftarrow \llbracket \langle \ell = \overrightarrow{v} \rangle \rrbracket_{J',\ell_i.p,\vec{\ell}}^*, I_{out} = I \rrbracket \\
&= \text{(Definition of } \llbracket - \rrbracket^* \text{)} \\
&= \llbracket c @ I_{in} \mid \langle I_{out}, c @ I_{in} \rangle @ I_{in} \leftarrow \llbracket v_i \rrbracket_{J',p,\vec{\ell}}^*, I_{out} = I \rrbracket \\
&= \text{(Definition of } \llbracket - \rrbracket^* \text{)} \\
&= \llbracket \llbracket v_i \rrbracket_{J',p,\downarrow,\vec{\ell},c} @ I_{in} \mid \langle I_{out}, c \rangle @ I_{in} \leftarrow \llbracket v_i \rrbracket_{J',p,\vec{\ell}}^*, I_{out} = I \rrbracket \\
&= \text{(Induction hypothesis)} \\
&= \llbracket v' . \vec{\ell} @ I_{in} \mid v' @ I_{in} \leftarrow \llbracket v_i \rrbracket_{J',p,I} \rrbracket \\
&= \text{(Definition of } \llbracket - \rrbracket^* \text{)} \\
&= \llbracket v' . \vec{\ell} @ I_{in} \mid v' @ I_{in} \leftarrow \llbracket \langle \ell = \overrightarrow{v} \rangle \rrbracket_{J',\ell_i.p,I} \rrbracket
\end{aligned}$$

Subcase $\downarrow.p$:

$$\begin{aligned}
& \llbracket \llbracket [v_i @ J_i]_{i=1}^n \rrbracket_{J',\downarrow,p,\downarrow,\vec{\ell},c} @ I_{in} \mid \langle I_{out}, c \rangle @ I_{in} \leftarrow \llbracket [v_i @ J_i]_{i=1}^n \rrbracket_{J',\downarrow,p,\vec{\ell}}^*, I_{out} = I \rrbracket \\
&= \text{(Definitions of } \llbracket - \rrbracket^* \text{ and } \llbracket - \rrbracket^* \text{)} \\
&= \text{concat}(\llbracket c @ I_{in} \mid \langle I_{out}, c @ I_{in} \rangle @ I_{in} \leftarrow \llbracket v_i \rrbracket_{J_i,p,\vec{\ell}}^*, I_{out} = I \rrbracket_{i=1}^n) \\
&= \text{(Induction hypothesis)} \\
&= \text{concat}(\llbracket v' . \vec{\ell} @ I_{in} \mid v' @ I_{in} \leftarrow \llbracket v_i \rrbracket_{J_i,p,I} \rrbracket_{i=1}^n) \\
&= \text{(Definition of } \llbracket - \rrbracket^* \text{)} \\
&= \llbracket v' . \vec{\ell} @ I_{in} \mid v' @ I_{in} \leftarrow \llbracket [v @ J]_{i=1}^n \rrbracket_{J',p} \rrbracket
\end{aligned}$$

Case $\langle \rangle$: The proof is the same as for base types with the constant c replaced by $\langle \rangle$.

Case $\langle \ell : \overrightarrow{A} \rangle$ where $|\vec{\ell}| \geq q$: We rely on the functions $\text{zip}_{\vec{\ell}}$, for transforming a record of lists of equal length to a list of records, and $\text{unzip}_{\vec{\ell}}$, for transforming a list of records to a record of lists of equal length. In fact we require special versions of zip and unzip that handle annotations, such that zip takes a record of lists of equal length whose annotations must be in sync, and unzip returns such a record.

$$\begin{aligned}
& \text{zip}_{\vec{\ell}}(\ell_i = [v_i]_{i=1}^n) = [v_i]_{i=1}^n \\
& \text{zip}_{\vec{\ell}}(\ell_i = v_i @ J :: s_i)_{i=1}^n = \langle \ell_i = v_i \rangle_{i=1}^n @ J :: \text{zip}_{\vec{\ell}}(\ell_i = s_i)_{i=1}^n \\
& \text{unzip}_{\vec{\ell}}(s) = \langle \ell_i = [v_i \cdot \ell_i @ J \mid v @ J \leftarrow s] \rangle_{i=1}^n
\end{aligned}$$

If $\vec{\ell}$ is a non-empty list of column labels then $\text{zip}_{\vec{\ell}}$ is the inverse of $\text{unzip}_{\vec{\ell}}$.

$$\text{zip}_{\vec{\ell}}(\text{unzip}_{\vec{\ell}}(s)) = s, \quad \text{if } |\vec{\ell}| \geq 1$$

$$\begin{aligned}
& \llbracket v \setminus_{J', p, \downarrow, \vec{\ell}, w} @I_{in} \mid \langle I_{out}, w \rangle @I_{in} \leftarrow \llbracket v \rrbracket_{J', p, \vec{\ell}, I_{out} = I}^* \\
= & \text{(Definition of } \setminus - \setminus \text{)} \\
& \langle \ell_i = \setminus v \setminus_{J', p, \downarrow, \vec{\ell}, \ell_i, w} \rangle_{i=1}^n @I_{in} \mid \langle I_{out}, w @I_{in} \rangle \leftarrow \llbracket v \rrbracket_{J', p, \vec{\ell}, \ell, I_{out} = I}^* \\
= & \text{(Definition of } zip \text{)} \\
& zip_{\vec{\ell}} \langle \ell_i = [v_{\ell_i, w} @I_{in} \mid \langle I_{out}, w @I_{in} \rangle \leftarrow \llbracket v \rrbracket_{J', p, \vec{\ell}, \ell_i, I_{out} = I}^*] \rangle_{i=1}^n \\
& \text{where } v_{\ell, w} \text{ stands for } \setminus v \setminus_{J', p, \downarrow, \vec{\ell}, \ell, w} \\
= & \text{(Induction hypothesis)} \\
& zip_{\vec{\ell}} \langle \ell_i = [v' \cdot \vec{\ell} \cdot \ell_i @I_{in} \mid v' @I_{in} \leftarrow \setminus v \setminus_{J', p, I}] \rangle_{i=1}^n \\
= & \text{(Definition of } zip \text{)} \\
& [v' \cdot \vec{\ell} @I_{in} \mid v' @I_{in} \leftarrow \setminus v \setminus_{J', p, I}]
\end{aligned}$$

Case Bag \vec{A} : Subcase ϵ : If $J' \neq I$ then both sides of the equation are equivalent to the empty bag. Suppose $J' = I$.

$$\begin{aligned}
& \llbracket [v_i @J_i]_{i=1}^n \setminus_{J', \downarrow, \vec{\ell}, I_{in}} @I_{in} \mid \langle I_{out}, I_{in} \rangle @I_{in} \leftarrow \llbracket [v_i @J_i]_{i=1}^n \rrbracket_{J', \vec{\ell}, I_{out} = I}^* \\
= & \text{(Definition of } \setminus - \setminus \text{)} \\
& [concat(\llbracket v_i \setminus_{J_i, \vec{\ell}, I_{in}} \rangle_{i=1}^n @I_{in} \mid \langle I_{out}, I_{in} \rangle @I_{in} \leftarrow \llbracket [v_i @J_i]_{i=1}^n \rrbracket_{J', \vec{\ell}, I_{out} = I}^*) \\
= & (\llbracket [v_i @J_i]_{i=1}^n \rrbracket_{J', \vec{\ell}} = \langle J', J_1 \rangle @J_1, \dots, \langle J', J_n \rangle @J_n) \\
& [concat(\llbracket v_i \setminus_{J_i, \vec{\ell}, I_{in}} \rangle_{i=1}^n @I_{in} \mid I_{in} \leftarrow \vec{J}, J' = I) \\
= & (J' = I) \\
& [concat(\llbracket v_i \setminus_{J_i, \vec{\ell}, I_{in}} \rangle_{i=1}^n @I_{in} \mid I_{in} \leftarrow \vec{J}) \\
= & \text{(Definition of } \setminus - \setminus \text{)} \\
& [concat(\llbracket v_i \setminus_{J_i, \vec{\ell}, I_{in}} \rangle_{i=1}^n @I_{in} \mid v_i @J_i \leftarrow [v_i @J_i]_{i=1}^n) @I_{in} \mid I_{in} \leftarrow \vec{J}) \\
= & \text{(Definition of } \setminus - \setminus \text{)} \\
& [concat(\llbracket v_i \cdot \vec{\ell} \setminus_{J_i, \epsilon, I_{in}} \mid v_i @J_i \leftarrow [v_i @J_i]_{i=1}^n) @I_{in} \mid I_{in} \leftarrow \vec{J}) \\
= & \text{(Definition of } \setminus - \setminus \text{)} \\
& [concat([v_i \cdot \vec{\ell} \mid v_i @J_i \leftarrow [v_i @J_i]_{i=1}^n, J_i = I_{in}) @I_{in} \mid I_{in} \leftarrow \vec{J}) \\
= & \text{(} v \text{ is well-indexed)} \\
& [v_i \cdot \vec{\ell} @I_{in} \mid v_i @I_{in} \leftarrow [v_i @J_i]_{i=1}^n] \\
= & \text{(Definition of } \setminus - \setminus \text{ and } J' = I) \\
& [v_i \cdot \vec{\ell} @I_{in} \mid v_i @I_{in} \leftarrow \llbracket [v_i @J_i]_{i=1}^n \setminus_{J', \epsilon, I} \rrbracket
\end{aligned}$$

Subcase ℓ_i, p :

$$\begin{aligned}
& \llbracket \langle \ell_i = v_i \rangle_{i=1}^n \setminus_{J', \ell_i, p, \downarrow, \vec{\ell}, I_{in}} @I_{in} \mid \langle I_{out}, I_{in} \rangle @I_{in} \leftarrow \llbracket \langle \ell_i = v_i \rangle_{i=1}^n \rrbracket_{J', \ell_i, p, \vec{\ell}, I_{out} = I}^* \\
= & \text{(Definitions of } \setminus - \setminus \text{ and } \llbracket - \rrbracket^* \text{)} \\
& \llbracket v_i \setminus_{J', p, \downarrow, \vec{\ell}, I_{in}} \rangle_{i=1}^n @I_{in} \mid \langle I_{out}, I_{in} \rangle @I_{in} \leftarrow \llbracket v_i \rrbracket_{J', p, \vec{\ell}, I_{out} = I}^* \\
= & \text{(Induction hypothesis)} \\
& [v' \cdot \vec{\ell} @I_{in} \mid v' @I_{in} \leftarrow \setminus v_i \setminus_{J', p, I}] \\
= & \text{(Definition of } \setminus - \setminus \text{)} \\
& [v' \cdot \vec{\ell} @I_{in} \mid v' @I_{in} \leftarrow \llbracket \langle \ell_i = v_i \rangle_{i=1}^n \setminus_{J', \ell_i, p, I} \rrbracket
\end{aligned}$$

Subcase \downarrow, p :

$$\begin{aligned}
& \llbracket [v_i @J_i]_{i=1}^n \setminus_{J', \downarrow, p, \downarrow, \vec{\ell}, I_{in}} @I_{in} \mid \langle I_{out}, I_{in} \rangle @I_{in} \leftarrow \llbracket [v_i @J_i]_{i=1}^n \rrbracket_{I_{out}, \downarrow, p, \vec{\ell}, I_{out} = I}^* \\
= & \text{(Definitions of } \setminus - \setminus \text{ and } \llbracket - \rrbracket^* \text{)} \\
& [(concat(\llbracket v_i \setminus_{J_i, p, \downarrow, \vec{\ell}, I_{in}} \rangle_{i=1}^n) @I_{in} \mid \langle I_{out}, I_{in} \rangle @I_{in} \leftarrow concat(\llbracket v_i \rrbracket_{J_i, p, \vec{\ell}}^*]_{i=1}^n), I_{out} = I) \\
= & \text{(Expanding } concat(\llbracket - \rrbracket^*]_{i=1}^n) \text{)} \\
& [(concat(\llbracket v_i \setminus_{J_i, p, \downarrow, \vec{\ell}, I_{in}} \rangle_{i=1}^n) @I_{in} \mid v_i @J_i \leftarrow [v_i @J_i]_{i=1}^n, \langle I_{out}, I_{in} \rangle @I_{in} \leftarrow \llbracket v_i \rrbracket_{J_i, p, \vec{\ell}}^*, I_{out} = I) \\
= & \text{(} v \text{ is well-indexed)} \\
& \llbracket v_i \setminus_{J_i, p, \downarrow, \vec{\ell}, I_{in}} \rangle_{i=1}^n @I_{in} \mid v_i @J_i \leftarrow [v_i @J_i]_{i=1}^n, \langle I_{out}, I_{in} \rangle @I_{in} \leftarrow \llbracket v_i \rrbracket_{J_i, p, \vec{\ell}}^*, I_{out} = I] \\
= & \text{(Comprehension } \rightarrow \text{ concatenation)} \\
& concat(\llbracket v_i \setminus_{J_i, p, \downarrow, \vec{\ell}, I_{in}} \rangle_{i=1}^n @I_{in} \mid \langle I_{out}, I_{in} \rangle @I_{in} \leftarrow \llbracket v_i \rrbracket_{J_i, p, \vec{\ell}}^*, I_{out} = I]_{i=1}^n) \\
= & \text{(Induction hypothesis)} \\
& concat([v' \cdot \vec{\ell} @I_{in} \mid v' @I_{in} \leftarrow \setminus v_i \setminus_{J_i, \downarrow, p, I}]_{i=1}^n) \\
= & \text{(Concatenation } \rightarrow \text{ comprehension)} \\
& [v' \cdot \vec{\ell} @I_{in} \mid v' @I_{in} \leftarrow \llbracket [v_i @J_i]_{i=1}^n \setminus_{J', \downarrow, p, I} \rrbracket
\end{aligned}$$

□

The proof of this lemma is the only part of the formalisation that makes use of values being well-indexed. Stitching shredded results together does not depend on any other property of indexes, thus any representation of indexes that yields unique indexes suffices.

Theorem 28. *If s is well-indexed and $\vdash \setminus s \setminus_{p, w} : A$ then $stitch_w(shred_{s, p}(A)) = \setminus s \setminus_{p, w}$.*

Proof. By induction on the structure of A .

Case O :

$$\begin{aligned}
& stitch_c(shred_{s, p}(O)) \\
= & c \\
= & \setminus s \setminus_{p, c}
\end{aligned}$$

Case $\langle \vec{\ell} = \vec{A} \rangle$:

$$\begin{aligned}
& stitch_{\langle \ell_i = w_i \rangle_{i=1}^n}(shred_{s, p}(\langle \ell_i : A \rangle_{i=1}^n)) \\
= & \langle \ell_i = stitch_{w_i}(shred_{s, p, i}(A)) \rangle_{i=1}^n \\
= & \text{(Induction hypothesis)} \\
& \langle \ell_i = \setminus s \setminus_{p, i, w_i} \rangle_{i=1}^n \\
= & \setminus s \setminus_{p, \langle \ell_i = w_i \rangle_{i=1}^n}
\end{aligned}$$

Case Bag A:

$$\begin{aligned}
& \text{stitch}_I(\text{shred}_{s,p}(\text{Bag } A)) \\
= & \text{stitch}_I(((\text{shred}_{s,p,\downarrow}(A))^{\llbracket s \rrbracket_p})) \\
= & [(\text{stitch}_n(\text{shred}_{s,p,\downarrow}(A)) @ I_{\text{in}}) \\
& \quad | \langle I_{\text{out}}, w \rangle @ I_{\text{in}} \leftarrow \llbracket s \rrbracket_p^*, I_{\text{out}} = I] \\
= & \quad (\text{Induction hypothesis}) \\
& \llbracket s \rrbracket_{p,\downarrow,w} @ I_{\text{in}} | \langle I_{\text{out}}, w \rangle @ I_{\text{in}} \leftarrow \llbracket s \rrbracket_p^*, I_{\text{out}} = I] \\
= & \quad (\text{Lemma 27}) \\
& \llbracket s \rrbracket_{p,I}
\end{aligned}$$

□

Proof of Theorem 7. By Theorem 28, setting $p = \epsilon$ and $n = \top \diamond 1$. □

We now obtain the main result.

Proof of Theorem 8. Recall the statement of Theorem 8: we need to show that $\vdash L : \text{Bag } A$ then:

$$\text{stitch}(\mathcal{H}[\llbracket L \rrbracket_{\text{Bag } A}]) = \mathcal{A}[\llbracket L \rrbracket]$$

and in particular:

$$\text{erase}(\text{stitch}(\mathcal{H}[\llbracket L \rrbracket_{\text{Bag } A}])) = \mathcal{N}[\llbracket L \rrbracket]$$

This follows immediately from Theorem 5, Lemma 6 and Theorem 7. □

D. Record flattening

Flat types

$$\text{Types } A, B ::= \text{Bag } \langle \vec{\ell} : \vec{O} \rangle$$

Flat terms

Query terms	$L, M ::= \text{let } \vec{C}$
Comprehensions	$C ::= \text{let } q = S \text{ in } S'$
Subqueries	$S ::= \text{for } (\vec{G} \text{ where } X) \text{ return } R$
Data sources	$u ::= t \mid q$
Generators	$G ::= x \leftarrow u$
Inner terms	$N ::= X \mid \text{index}$
Record terms	$R ::= \langle \vec{\ell} = \vec{N} \rangle$
Base terms	$X ::= x.\ell \mid c(\vec{X})$

Flattening types

$$(\text{Bag } A)^\succ = \text{Bag } (A^\succ)$$

$$\begin{aligned}
O^\succ &= \langle \bullet : O \rangle \\
\langle \langle \vec{\ell} : \vec{F} \rangle \rangle^\succ &= \langle \{ \langle \ell_i \ell'_i \rangle : O \mid i \leftarrow \text{dom}(\vec{\ell}), (\ell'_i : O) \leftarrow F_i^\succ \} \rangle
\end{aligned}$$

Flattening terms

$$\begin{aligned}
(\text{let } q = S \text{ in } S')^\succ &= \text{let } q = S \text{ in } S'^\succ \\
(\text{for } (\vec{G} \text{ where } X) \text{ return } N)^\succ &= \text{for } (\vec{G} \text{ where } X) \text{ return } N^\succ \\
X^\succ &= \langle \bullet = X \rangle \\
\langle \langle \ell_i = N_i \rangle_{i=1}^m \rangle^\succ &= \langle \langle \ell_i \ell'_j = X_j \rangle_{(i=1,j=1)}^{(m,n_i)} \rangle, \\
&\quad \text{where } N_i^\succ = \langle \ell'_j = X_j \rangle_{j=1}^{n_i}
\end{aligned}$$

Type soundness

$$\vdash L : A \Rightarrow \vdash L^\succ : A^\succ$$

Flat values

$$\begin{aligned}
\text{Results } s &::= \vec{r} \\
\text{Rows } r &::= \langle \ell_1 = c_1, \dots, \ell_n = c_n \rangle
\end{aligned}$$

Semantics of flat record queries

$$\mathcal{R}[\llbracket L \rrbracket] = \mathcal{L}[\llbracket L \rrbracket]$$

Unflattening record values

$$[r_1, \dots, r_n]^\prec = [(r_1)^\prec, \dots, (r_n)^\prec]$$

$$\begin{aligned}
\langle \bullet = c \rangle^\prec &= c \\
\langle \langle \ell_i \ell'_j = c_j \rangle_{(i=1,j=1)}^{(m,n_i)} \rangle^\prec &= \langle \ell_i = \langle \langle \ell'_j = c_j \rangle_{j=1}^{n_i} \rangle_{i=1}^m \rangle
\end{aligned}$$

Type soundness

$$\vdash L : A \Rightarrow \vdash L^\prec : A^\prec$$

Correctness

Proposition 29. If L is a let-inserted query and $\vdash L : \text{Bag } A$, then

$$(\mathcal{R}[\llbracket L^\succ \rrbracket])^\prec = \mathcal{L}[\llbracket L \rrbracket]$$