# Learning Models over Relational Databases

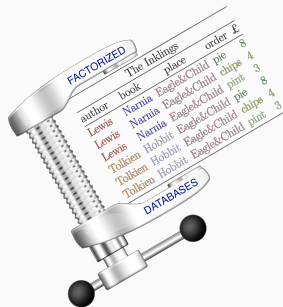fdbresearch.github.io                    relational.ai

---

**Dan Olteanu**

Oxford & relationalAI

FG DB Symposium 2020
Darmstadt, March 2020

# Acknowledgments

**Kaggle Survey:** Most Data Scientists use Relational Data at Work!



Overall

By Industry

# Relational Model: Jewel in the Data Management Crown

- Massive adoption of the Relational Model in last decades

- Many human hours invested in building relational models

- Relational databases are rich with knowledge of the underlying domains

# Current State of Affairs in Analytics Workloads



- Carefully crafted by domain experts
- Comes with relational structure

- Throws away relational structure
- Can be order-of-magnitude larger

# Conjecture

The learning time and accuracy of the model can be drastically improved by exploiting the structure and semantics of the underlying multi-relational database.

# Current Landscape for ML over DB

## No integration



**The good:**

1. Most DB+ML solutions operate in this space
2. Supports virtually any ML task
3. ML & DB distinct tools on the technology stack

**The bad:**

1. Materialisation of feature extraction query
2. DB exports data as one table, ML imports it in own format
3. One/multi-hot encoding of categorical variables

**Examples:**

PostgreSQL + R,     Pandas + scikit-learn/TensorFlow,     SparkSQL + MLlib, etc.

Database System

Feature Extraction Query → Query Eval → materialised output = data matrix → Model Learning → $\theta$

Model

- DB supports ML tasks as UDF
- Same running process for DB and ML
- DB computes one table, ML works directly on it → No data export/import

**Examples:**

MadLib supports comprehensive library of ML UDFs

Bismark gives unified programming architecture for incremental gradient descent

**Structure-Aware Learning** vs. **Structure-Agnostic Learning**

- Exploit relational structure and semantics
- Exploit database optimisations, e.g., push parts of ML tasks past joins
- One evaluation plan for mixed DB and ML workload

**Structure-aware** Learning **FASTER** even than Feature Extraction Query!

## Case in Point (1): A Retailer Use Case



| Relation | Cardinality | Arity (Keys+Values) | File Size (CSV) |
|---|---|---|---|
| Inventory | 84,055,817 | 3 + 1 | 2 GB |
| Items | 5,618 | 1 + 4 | 129 KB |
| Stores | 1,317 | 1 + 14 | 139 KB |
| Demographics | 1,302 | 1 + 15 | 161 KB |
| Weather | 1,159,457 | 2 + 6 | 33 MB |
| Join | 84,055,817 | 3 + 41 | 23GB |

Train a linear regression model to predict *inventory* given all features

|  | PostgreSQL+TensorFlow | |
|  | **Time** | **Size (CSV)** |
| --- | --- | --- |
| Database | – | 2.1 GB |
| Join | 152.06 secs | 23 GB |
| Export | 351.76 secs | 23 GB |
| Shuffling | 5,488.73 secs | 23 GB |
| Query batch | – | – |
| Grad Descent | 7,249.58 secs | – |
| Total time | 13,242.13 secs | |

Train a linear regression model to predict *inventory* given all features

| | PostgreSQL+TensorFlow | | Our approach (SIGMOD'19) | |
| | **Time** | **Size (CSV)** | **Time** | **Size (CSV)** |
| Database | – | 2.1 GB | – | 2.1 GB |
| Join | 152.06 secs | 23 GB | – | – |
| Export | 351.76 secs | 23 GB | – | – |
| Shuffling | 5,488.73 secs | 23 GB | – | – |
| Query batch | – | – | 6.08 secs | 37 KB |
| Grad Descent | 7,249.58 secs | – | 0.05 secs | – |
| Total time | 13,242.13 secs | | 6.13 secs | |

$2,160\times$ faster while being more accurate (RMSE on 2% test data)

Train a linear regression model to predict *inventory* given all features

| | PostgreSQL+TensorFlow | | Our approach (SIGMOD'19) | |
| | Time | Size (CSV) | Time | Size (CSV) |
| --- | --- | --- | --- | --- |
| Database | – | 2.1 GB | – | 2.1 GB |
| Join | 152.06 secs | 23 GB | – | – |
| Export | 351.76 secs | 23 GB | – | – |
| Shuffling | 5,488.73 secs | 23 GB | – | – |
| Query batch | – | – | 6.08 secs | 37 KB |
| Grad Descent | 7,249.58 secs | – | 0.05 secs | – |
| Total time | 13,242.13 secs | | 6.13 secs | |

$2,160\times$ faster while being more accurate (RMSE on 2% test data)

TensorFlow trains one model. Our approach takes $< 0.1$ sec for any extra model over a subset of the given feature set.

**TensorFlow's Behaviour is the Rule, not the Exception!**

Similar behaviour (or outright failure) for more:

- **datasets**: Favorita, TPC-DS, Yelp, Housing

- **systems**:
    - used in industry: R, scikit-learn, Python StatsModels, mlpack, XGBoost, MADlib

    - academic prototypes: Morpheus, libFM

- **models**: decision trees, factorisation machines, $k$-means, ..

**This is to be contrasted with the scalability of DBMSs!**

# How to achieve this performance improvement?

# Idea 1: Turn the ML Problem into a DB Problem

## Through DB Glasses, Everything is a Batch of Queries

| Workload | Query Batch | # Queries |
|---|---|---|
| Linear Regression Covariance Matrix | $\text{SUM}(X_i * X_j)$ $\text{SUM}(X_i)$ $\text{GROUP BY}$ $X_j$ $\text{SUM(1)}$ $\text{GROUP BY}$ $X_i, X_j$ | 814 |
| Decision Tree (Regression, 1 Node) | $\text{VARIANCE}(Y)$ $\text{WHERE}$ $X_j = c_j$ | 3,141 |
| R$k$-means | $\text{SUM(1)}$ $\text{GROUP BY}$ $X_j$ $\text{SUM(1)}$ $\text{GROUP BY}$ $\text{Center}_1, \ldots, \text{Center}_k$ | 41 |

(# Queries shown for Retailer dataset with 39 attributes)

Queries in a batch:

- Same aggregates but over different attributes
- Expressed over the same join of the database relations

AMPLE opportunities for sharing computation in a batch.

## Models under Consideration

So far:

- Polynomial regression
- Factorisation machines
- Classification/regression trees
- Mutual information
- Chow Liu trees
- $k$-means clustering
- $k$-nearest neighbours
- (robust, ordinal) PCA
- SVM

On-going:

- Boosting regression trees
- AdaBoost
- Sum-product networks
- Random forests
- Logistic regression
- Linear algebra:
  - QR decomposition
  - SVD
  - low-rank matrix factorisation

All these cases can benefit from **structure-aware computation**

**Ridge Linear Regression**

$\Downarrow$

**Query Batch**

## Recap: Ridge Linear Regression

Linear regression model:

$$f_{\boldsymbol{\theta}}(\mathbf{x}) = \langle \boldsymbol{\theta}, \mathbf{x} \rangle = \theta_0 x_0 + \theta_1 x_1 + \ldots$$

- Training dataset $D$ defined by *feature extraction query*
  - A tuple $(\mathbf{x}, y) \in D$ consists of feature vector $\mathbf{x}$ and response $y$

- Parameters $\boldsymbol{\theta}$ obtained by minimising the objective function:

$$J(\boldsymbol{\theta}) = \overbrace{\frac{1}{2|D|} \sum_{(\mathbf{x},y) \in D} (\langle \boldsymbol{\theta}, \mathbf{x} \rangle - y)^2}^{\text{least square loss}} + \overbrace{\frac{\lambda}{2} \|\boldsymbol{\theta}\|_2^2}^{\ell_2 - \text{regulariser}}$$

We can solve $\theta^* := \arg\min_{\theta} J(\theta)$ with batch-gradient descent:

```
repeat until convergence:
```
$$\theta := \theta - \alpha \cdot \nabla J(\theta)$$

**Model reformulation idea**: Decouple

- data-dependent ($\mathbf{x}$, $y$) computation from
- data-independent ($\theta$) computation

in the formulations of the objective $J(\theta)$ and its gradient $\nabla J(\theta)$.

$$J(\boldsymbol{\theta}) = \frac{1}{2|D|} \sum_{(\mathbf{x},y) \in D} (\langle \boldsymbol{\theta}, \mathbf{x} \rangle - y)^2 + \frac{\lambda}{2} \|\boldsymbol{\theta}\|_2^2$$

$$= \frac{1}{2|D|} \left( \boldsymbol{\theta}^\top \underbrace{\left( \sum_{(\mathbf{x},y) \in D} \mathbf{x}\mathbf{x}^\top \right)}_{\Sigma} \boldsymbol{\theta} - 2 \Big\langle \boldsymbol{\theta}, \underbrace{\sum_{(\mathbf{x},y) \in D} y \cdot \mathbf{x}}_{\mathbf{c}} \Big\rangle + \underbrace{\left( \sum_{(\mathbf{x},y) \in D} y^2 \right)}_{s_Y} \right) + \frac{\lambda}{2} \|\boldsymbol{\theta}\|_2^2$$

$$= \frac{1}{2|D|} \left( \boldsymbol{\theta}^\top \Sigma \boldsymbol{\theta} - 2 \langle \boldsymbol{\theta}, \mathbf{c} \rangle + s_Y \right) + \frac{\lambda}{2} \|\boldsymbol{\theta}\|_2^2$$

$$\nabla J(\boldsymbol{\theta}) = \frac{1}{|D|} \left( \Sigma \boldsymbol{\theta} - \mathbf{c} \right) + \lambda \boldsymbol{\theta}$$

Compute one query for each entry $\sum_{(\mathbf{x},y) \in D} \mathbf{x}_i \mathbf{x}_j^\top$ in $\Sigma$:

Compute one query for each entry $\sum_{(\mathbf{x},y) \in D} \mathbf{x}_i \mathbf{x}_j^\top$ in $\Sigma$:

- $x_i$, $x_j$ continuous

  ```sql
  SELECT SUM (x_i * x_j) FROM D;
  ```

where $D$ is the feature extraction query over the input DB.

Compute one query for each entry $\sum_{(\mathbf{x},y) \in D} \mathbf{x}_i \mathbf{x}_j^\top$ in $\Sigma$:

- $x_i$, $x_j$ continuous

  ```sql
  SELECT SUM (x_i * x_j) FROM D;
  ```

- $\mathbf{x}_i$ categorical, $x_j$ continuous

  ```sql
  SELECT x_i, SUM(x_j) FROM D GROUP BY x_i;
  ```

where $D$ is the feature extraction query over the input DB.

Compute one query for each entry $\sum_{(\mathbf{x},y) \in D} \mathbf{x}_i \mathbf{x}_j^\top$ in $\Sigma$:

- $x_i$, $x_j$ continuous

  ```
  SELECT SUM (x_i * x_j) FROM D;
  ```

- $\mathbf{x}_i$ categorical, $x_j$ continuous

  ```
  SELECT x_i, SUM(x_j) FROM D GROUP BY x_i;
  ```

- $\mathbf{x}_i$, $\mathbf{x}_j$ categorical

  ```
  SELECT x_i, x_j, SUM(1) FROM D GROUP BY x_i, x_j;
  ```
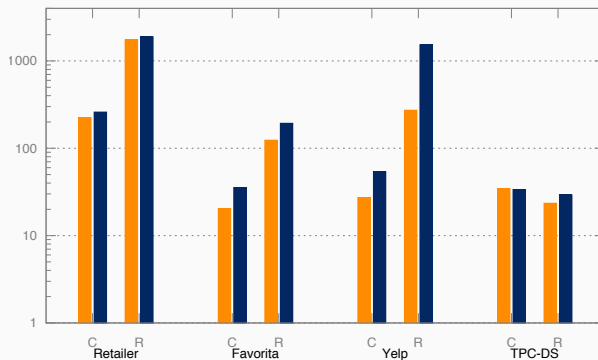
where $D$ is the feature extraction query over the input DB.

**Natural Attempt:**

Use Existing DB System to Compute Query Batch

# Existing DBMSs are NOT Designed for Query Batches



Relative Speedup for **Our Approach** over **DBX** and **MonetDB**

C = Covariance Matrix; R = Regression Tree Node; AWS d2.xlarge (4 vCPUs, 32GB)

# Idea 2: Exploit Problem Structure to Lower Complexity

Algebraic structure: (semi)rings $(\mathcal{R}, +, *, \mathbf{0}, \mathbf{1})$

- Distributivity law $\rightarrow$ Factorisation

  Factorised Databases                                    [VLDB'12+'13,TODS'15,SIGREC'16]

  Factorised Machine Learning      [SIGMOD'16+'19,DEEM'18,PODS'18+'19, TODS'20]

- Additive inverse $\rightarrow$ Uniform treatment of updates

  Factorised Incremental Maintenance                                    [SIGMOD'18+'20]

- Sum-Product abstraction $\rightarrow$ Same processing for distinct tasks

  DB queries, Covariance matrix, PGM inference, Matrix chain multiplication

                                                                        [SIGMOD'18+'19]

Combinatorial structure: query width and data degree measures

- Width measure $w$ for FEQ $\rightarrow$ Low complexity $\tilde{O}(N^w)$

  factorisation width $\geq$ fractional hypertree width $\geq$ sharp-submodular width
  worst-case optimal size and time for factorised joins

  [ICDT'12+'18,TODS'15,PODS'19,TODS'20]

- Degree $\rightarrow$ Adaptive processing depending on high/low degrees

  worst-case optimal incremental maintenance                [ICDT'19a, PODS'20]
  evaluation of queries with negated relations of bounded degree        [ICDT'19b]

- Functional dependencies $\rightarrow$ Learn simpler, equivalent models

  reparameterisation of polynomial regression models and factorisation machines

  [PODS'18,TODS'20]

# Idea 3: Lower the Constant Factors

1. Specialisation for workload and data

   Generate code specific to the query batch and dataset

   Improve cache locality for hot data path

2. Sharing low-level data access

   Aggregates decomposed into views over join tree
   Share data access across views with different output schemas

3. Parallelisation: multi-core (SIMD & distribution to come)

   Task and domain parallelism

   [DEEM'18,SIGMOD'19, CGO'20]

Code Optimisations

⇓

Non-trivial Speedup

## One DSL to Express both DB and ML Workloads! [CGO'20]

**Collections are Dictionaries or Sets**

- Database relations are modeled as dictionaries

| Relation R(A,B) | |
|---|---|
| A | B |
| $a_1$ | $b_1$ |
| $a_1$ | $b_1$ |
| $a_2$ | $b_1$ |
| $a_2$ | $b_1$ |
| $a_2$ | $b_2$ |

| Relation R(A,B) in IFAQ | | | |
|---|---|---|---|
| A | B | $\rightarrow$ | $R(A, B)$ |
| $a_1$ | $b_1$ | $\rightarrow$ | 2 |
| $a_2$ | $b_1$ | $\rightarrow$ | 2 |
| $a_2$ | $b_2$ | $\rightarrow$ | 1 |

Inspired by the FAQ framework [PODS'16]

## IFAQ: Iterative Functional Aggregate Queries

- $\Sigma$ for stateful computation over collection elements:

| **IFAQ** | | **C++** |
|---|---|---|
| $\displaystyle\sum_{e \in set} f(e)$ | $\xrightarrow{\text{Compile}}$ | `for(auto& e : set)`<br>`    res += f(e);` |

- $\lambda$ for constructing dictionaries:

| **IFAQ** | | **C++** |
|---|---|---|
| $\displaystyle\lambda_{e \in set} f(e)$ | $\xrightarrow{\text{Compile}}$ | `for(auto& e : set)`<br>`    res[e] = f(e);` |

- Supports while loops and conditionals

# Transformation Steps for IFAQ Expressions

**Dataset with three relations:**

    **S**ales(item,store,unit sales)       **I**tem(item, price)       Sto**R**e(store, city)

**Learning Task:**

      Learn Linear Regression model to predict number of **u**nit sales.

**Training Dataset:**

$$Q(x) = S(x_S) \bowtie R(x_R) \bowtie I(x_I)$$

> **Batch Gradient Descent:**
> Update $\theta$ in direction of gradient of square loss

```
let F = [[i, s, p, c]] in
```

$\theta \leftarrow \theta_0$

```
while( not converged ) {
```

$$\theta = \lambda_{f_1 \in \mathbf{F}} \left( \theta(f_1) - \frac{\alpha}{|\mathbf{Q}|} \underbrace{\sum_{x \in \mathrm{sup}(\mathbf{Q})} \mathbf{Q}(x) * \left( \sum_{f_2 \in \mathbf{F}} \theta(f_2) * x[f_2] - x[u] \right) * x[f_1]}_{\text{Gradient of square loss}} \right)$$

```
}
```

$\theta$

# (Simplified) Linear Regression in IFAQ

> **Batch Gradient Descent:**
> Update $\theta$ in direction of gradient of square loss

```
let F = [[i, s, p, c]] in
```
$$\theta \leftarrow \theta_0$$
```
while( not converged ) {
```

$$\theta = \lambda_{f_1 \in \mathbf{F}} \left( \theta(f_1) - \frac{\alpha}{|\mathbf{Q}|} \underbrace{\sum_{x \in \sup(\mathbf{Q})} \mathbf{Q}(x) * \left( \sum_{f_2 \in \mathbf{F}} \theta(f_2) * x[f_2] - x[u] \right) * x[f_1]}_{\text{Gradient of square loss}} \right)$$

```
}
```
$$\theta$$

> For simplicity and WLOG, we
> 1. set $\frac{\alpha}{|\mathbf{Q}|} = 1$
> 2. ignore $x[u]$

```
let F = [[i, s, p, c]] in
```

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta}_0$$

```
while( not converged ) {
```

$$\boldsymbol{\theta} = \lambda_{f_1 \in \mathbf{F}} \left( \boldsymbol{\theta}(f_1) - \sum_{x \in \sup(\mathbf{Q})} \mathbf{Q}(x) * \left( \sum_{f_2 \in \mathbf{F}} \boldsymbol{\theta}(f_2) * x[f_2] \right) * x[f_1] \right)$$

```
}
```

$$\boldsymbol{\theta}$$

Transformation Rule: **Normalisation**

$$\theta = \lambda_{f_1 \in \mathbf{F}} \left( \theta(f_1) \ - \ \sum_{x \in \mathrm{sup}(\mathbf{Q})} \mathbf{Q}(x) \ * \ \sum_{f_2 \in \mathbf{F}} \left( \ \theta(f_2) \ * \ x[f_2] \ \right) \ * \ x[f_1] \right)$$

Transformation Rule: **Normalisation**

$$\theta = \lambda_{f_1 \in \mathbf{F}} \left( \theta(f_1) \; - \sum_{x \in \sup(\mathbf{Q})} \mathbf{Q}(x) \; * \sum_{f_2 \in \mathbf{F}} \left( \; \theta(f_2) \; * \; x[f_2] \; \right) \; * \; x[f_1] \right)$$

Transformation Rule: **Normalisation**

$$\theta = \lambda_{f_1 \in \mathbf{F}} \left( \theta(f_1) \; - \sum_{x \in \mathrm{sup}(\mathbf{Q})} \; \sum_{f_2 \in \mathbf{F}} \left( \mathbf{Q}(x) \; * \; \theta(f_2) \; * \; x[f_2] \; * \; x[f_1] \right) \right)$$

Transformation Rule: **Loop Scheduling**

$$\theta = \underset{f_1 \in \mathbf{F}}{\lambda} \left( \theta(f_1) \; - \; \sum_{x \in \operatorname{sup}(\mathbf{Q})} \; \sum_{f_2 \in \mathbf{F}} \Big( \mathbf{Q}(x) \; * \; \theta(f_2) \; * \; x[f_2] \; * \; x[f_1] \Big) \right)$$

Order loops by size of support

Transformation Rule: **Loop Scheduling**

$$\theta = \lambda_{f_1 \in \mathbf{F}} \left( \theta(f_1) \; - \sum_{f_2 \in \mathbf{F}} \; \sum_{x \in \sup(\mathbf{Q})} \Big( \mathbf{Q}(x) \; * \; \theta(f_2) \; * \; x[f_2] \; * \; x[f_1] \Big) \right)$$

Order loops by size of support

Transformation Rule: **Factorisation**

$$\theta = \lambda_{f_1 \in \mathbf{F}} \left( \theta(f_1) \; - \sum_{f_2 \in \mathbf{F}} \; \sum_{x \in \sup(\mathbf{Q})} \left( \mathbf{Q}(x) \; * \; \theta(f_2) \; * \; x[f_2] \; * \; x[f_1] \right) \right)$$

Transformation Rule: **Factorisation**

$$\theta = \lambda_{f_1 \in \mathbf{F}} \left( \theta(f_1) \; - \sum_{f_2 \in \mathbf{F}} \; \sum_{x \in \mathrm{sup}(\mathbf{Q})} \left( \mathbf{Q}(x) \; * \; \theta(f_2) \; * \; x[f_2] \; * \; x[f_1] \right) \right)$$

Less arithmetic operations

Transformation Rule: **Factorisation**

$$\theta = \lambda_{f_1 \in \mathbf{F}} \left( \theta(f_1) \; - \sum_{f_2 \in \mathbf{F}} \; \theta(f_2) \; * \sum_{x \in \sup(\mathbf{Q})} \; \Big( \mathbf{Q}(x) \; * \; x[f_2] \; * \; x[f_1] \Big) \right)$$

Less arithmetic operations

Transformation Rule: **Static Memoisation**

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta}_0$$

```
while( not converged ){
```
$$\boldsymbol{\theta} = \lambda_{f_1 \in \mathbf{F}} \left( \boldsymbol{\theta}(f_1) - \sum_{f_2 \in \mathbf{F}} \boldsymbol{\theta}(f_2) * \sum_{x \in \mathrm{sup}(\mathbf{Q})} \left( \mathbf{Q}(x) * x[f_2] * x[f_1] \right) \right)$$
```
}
```
$$\boldsymbol{\theta}$$

Transformation Rule: **Static Memoisation**

$$\theta \leftarrow \theta_0$$

```
while( not converged ){
```

$$\theta = \lambda_{f_1 \in \mathbf{F}} \left( \theta(f_1) - \sum_{f_2 \in \mathbf{F}} \theta(f_2) * \sum_{x \in \sup(\mathbf{Q})} \Big( \mathbf{Q}(x) * x[f_2] * x[f_1] \Big) \right)$$

```
}
```

$$\theta$$

Transformation Rule: **Code Motion**

$$\theta \leftarrow \theta_0$$

```
while( not converged ){
```

$$\text{let } \mathbf{M} = \lambda_{f_1 \in \mathbf{F}} \lambda_{f_2 \in \mathbf{F}} \sum_{x \in \text{sup}(\mathbf{Q})} \mathbf{Q}(x) * x[f_2] * x[f_1] \text{ in}$$

$$\theta = \lambda_{f_1 \in \mathbf{F}} \Big( \theta(f_1) - \sum_{f_2 \in \mathbf{F}} \theta(f_2) * \mathbf{M}(f_1)(f_2) \Big)$$

```
}
```

$$\theta$$

**M** defines the covariance matrix

Transformation Rule: **Code Motion**

$$\theta \leftarrow \theta_0$$

```
while( not converged ){
```

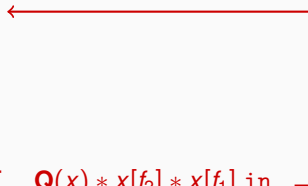$$\text{let } \mathbf{M} = \lambda_{f_1 \in \mathbf{F}} \lambda_{f_2 \in \mathbf{F}} \sum_{x \in \text{sup}(\mathbf{Q})} \mathbf{Q}(x) * x[f_2] * x[f_1] \text{ in}$$

$$\theta = \lambda_{f_1 \in \mathbf{F}} \left( \theta(f_1) - \sum_{f_2 \in \mathbf{F}} \theta(f_2) * \mathbf{M}(f_1)(f_2) \right)$$

```
}
```

$$\theta$$

Transformation Rule: **Code Motion**

```
let M =
```
$$\lambda_{f_1 \in \mathbf{F}} \lambda_{f_2 \in \mathbf{F}} \sum_{x \in \mathrm{sup}(\mathbf{Q})} \mathbf{Q}(x) * x[f_1] * x[f_2]$$ `in`

$$\theta \leftarrow \theta_0$$

```
while( not converged ){
```

$$\theta = \lambda_{f_1 \in \mathbf{F}} \left( \theta(f_1) - \sum_{f_2 \in \mathbf{F}} \theta(f_2) * \mathbf{M}(f_1)(f_2) \right)$$

```
}
```

$$\theta$$

Expression after High-Level Optimisations:

$$\texttt{let } \mathbf{M} = \lambda_{f_1 \in \mathbf{F}} \lambda_{f_2 \in \mathbf{F}} \sum_{x \in \text{sup}(\mathbf{Q})} \mathbf{Q}(x) * x[f_1] * x[f_2] \texttt{ in}$$

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta}_0$$

$$\texttt{while( not converged )} \{$$

$$\boldsymbol{\theta} = \lambda_{f_1 \in \mathbf{F}} \left( \boldsymbol{\theta}(f_1) - \sum_{f_2 \in \mathbf{F}} \boldsymbol{\theta}(f_2) * \mathbf{M}(f_1)(f_2) \right)$$

$$\}$$

$$\boldsymbol{\theta}$$

Transformation Rule: **Loop Unrolling**

```
let M = λ   λ   Σ      Q(x) * x[f₁] * x[f₂] in
       f₁∈F f₂∈F x∈sup(Q)
```

$$\theta \leftarrow \theta_0$$

```
while( not converged ){
```

$$\theta = \underset{f_1 \in \mathbf{F}}{\lambda} \left( \theta(f_1) - \sum_{f_2 \in \mathbf{F}} \theta(f_2) * \mathbf{M}(f_1)(f_2) \right)$$

```
}
```

$$\theta$$

Unroll Loops over statically known features **F**

Transformation Rule: **Loop Unrolling**

```
let M =
```
$$\lambda_{f_1 \in \mathbf{F}} \lambda_{f_2 \in \mathbf{F}} \sum_{x \in \mathrm{sup}(\mathbf{Q})} \mathbf{Q}(x) * x[f_1] * x[f_2] \texttt{ in}$$
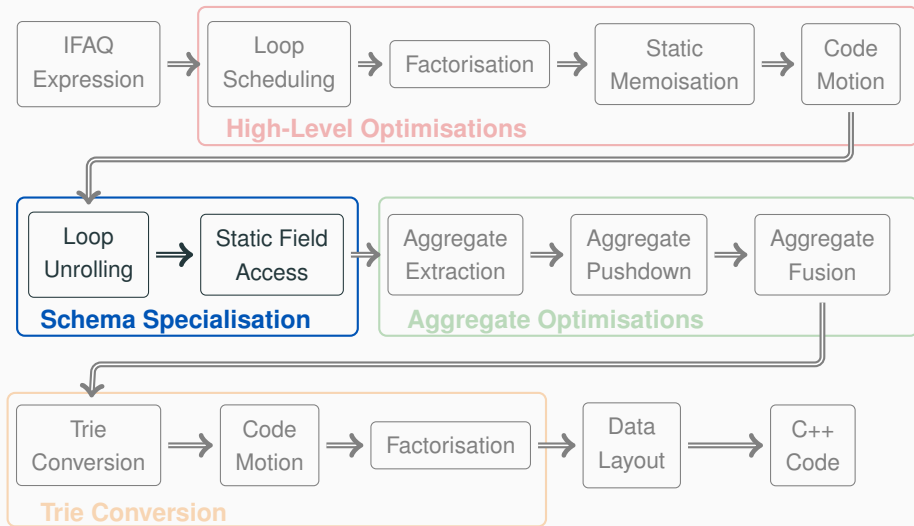
$$\theta \leftarrow \theta_0$$

```
while( not converged ){
```

$$\theta = \left\{ \left\{ c \rightarrow \left( \theta(c) - (... + \theta(c) * \mathbf{M}(c)(c) + \theta(p) * \mathbf{M}(c)(p)...) \right), ... \right\} \right\}$$

```
}
```

$$\theta$$

> Unroll Loops over statically known features **F**

Transformation Rule: **Loop Unrolling**

$$\texttt{let } \mathbf{M} = \left\{ \left\{ c \rightarrow \left\{ \left\{ ..., p \rightarrow \sum_{x \in \sup(\mathbf{Q})} \mathbf{Q}(x) * x[c] * x[p], ... \right\} \right\}, ... \right\} \right\} \texttt{ in}$$

$$\theta \leftarrow \theta_0$$

$$\texttt{while( not converged )\{}$$

$$\theta = \left\{ \left\{ c \rightarrow \left( \theta(c) - \left( ... + \theta(c) * \mathbf{M}(c)(c) + \theta(p) * \mathbf{M}(c)(p)... \right) \right), ... \right\} \right\}$$

$$\texttt{\}}$$

$$\theta$$

- Convert dictionaries over **F** into records
- Dynamic accesses into static accesses

Transformation Rule: **Static Field Access**

$$\texttt{let } M = \Big\{ c = \big\{ ..., p = \sum_{x \in \sup(\mathbf{Q})} \mathbf{Q}(x) * x.c * x.p, ..., \big\}, ... \big\} \texttt{ in}$$

$$\theta \leftarrow \theta_0$$

$$\texttt{while( not converged )\{}$$

$$\theta = \Big\{ c = \theta.c - \Big( ... + \theta.c * M.c.c + \theta.p * M.c.p... \Big), ... \Big\}$$

$$\texttt{\}}$$

$$\theta$$

- Convert dictionaries over **F** into records
- Dynamic accesses into static accesses

Transformation Rule: **Aggregate Extraction**

```
let M =
    {c =
        {..., c = ∑       Q(x) ∗ x.c ∗ x.c,  p = ∑       Q(x) ∗ x.c ∗ x.p, ...}
              x∈sup(Q)                            x∈sup(Q)
    , ...} in ...
```

Transformation Rule: **Aggregate Extraction**

```
let M =
    { c =
```

$$\left\{ ..., c = \sum_{x \in \mathrm{sup}(\mathbf{Q})} \mathbf{Q}(x) * x.c * x.c, \ p = \sum_{x \in \mathrm{sup}(\mathbf{Q})} \mathbf{Q}(x) * x.c * x.p, ... \right\}$$

```
    , ... } in ...
```

Transformation Rule: **Aggregate Extraction**

```
let M_cc =  ∑       Q(x) * x.c * x.c in
          x∈sup(Q)

let M_cp =  ∑       Q(x) * x.c * x.p in
          x∈sup(Q)

let M = {    c = {...,   c = M_cc,   p = M_cp,   ...},...} in ...
```

Recall: $\mathbf{Q}(x) = S(x_S) \bowtie I(x_I) \bowtie R(x_R)$

Transformation Rule: **Aggregate Extraction**

```
let M_cc =  ∑      Q(x) * x.c * x.c in
           x∈sup(Q)
```

```
let M_cp =  ∑      Q(x) * x.c * x.p in
           x∈sup(Q)
```

```
let M = {     c = {...,   c = M_cc,   p = M_cp,   ...},...} in ...
```

Recall: $\mathbf{Q}(x) = S(x_S) \bowtie I(x_I) \bowtie R(x_R)$

Transformation Rule: **Aggregate Extraction**

```
let M_cc =   ∑    Q(x) * x.c * x.c in
           x∈sup(Q)
```

$$\texttt{let } M_{cc} = \sum_{x \in \mathrm{sup}(\mathbf{Q})} \mathbf{Q}(x) * x.c * x.c \texttt{ in}$$

$$\texttt{let } M_{cp} = \sum_{x \in \mathrm{sup}(\mathbf{Q})} \mathbf{Q}(x) * x.c * x.p \texttt{ in}$$

$$\texttt{let } M = \{ \quad c = \{..., \quad c = M_{cc}, \quad p = M_{cp}, \quad ...\}, ...\} \texttt{ in } ...$$

We can:

- avoid materialisation of $\mathbf{Q}(x)$
- inline code for join computation

To compute $\mathbf{S}(x_S) \bowtie \mathbf{R}(x_R)$ on variable $s$:

1. Construct nested dictionaries over $\mathbf{R}$:

$$\mathbf{H}_R = \sum_{x_R \in \sup(\mathbf{R})} \mathbf{R}(x_R) * \{\{\{s = x_R.s\} \rightarrow \{\{x_R \rightarrow 1\}\}\}\}$$

For join value $s$, $\mathbf{H}_R(s)$ maps to partition of $\mathbf{R}$ with $s = x_R.s$

To compute $\mathbf{S}(x_S) \bowtie \mathbf{R}(x_R)$ on variable $s$:

1. Construct nested dictionaries over $\mathbf{R}$:

$$\mathbf{H}_R = \sum_{x_R \in \sup(\mathbf{R})} \mathbf{R}(x_R) * \{\{\{s = x_R.s\} \rightarrow \{\{x_R \rightarrow 1\}\}\}\}$$

For join value $s$, $\mathbf{H}_R(s)$ maps to partition of $\mathbf{R}$ with $s = x_R.s$

2. Iterate over $\mathbf{S}$, and probe $\mathbf{H}_R$ for joining tuples:

$$\mathbf{J}_{S \bowtie R} = \sum_{x_S \in \sup(\mathbf{S})} \sum_{x_R \in \sup(\mathbf{H}_R(\{s = x_S.s\}))}$$

$$\text{let } k = \{s = x_S.s, \ i = x_S.i, \ u = x_S.u, \ c = x_R.c\} \text{ in}$$

$$\{\{k \rightarrow \mathbf{S}(x_S) * \mathbf{H}_R(\{s = x_S.s\})(x_R)\}\}$$

Transformation Rule: **Aggregate Pushdown**

$$\mathbf{H}_R = \sum_{x_R \in \sup(\mathbf{R})} \mathbf{R}(x_R) * \{\{\{s = x_R.s\} \rightarrow \{\{x_R \rightarrow 1\}\}\}\}$$

$$\mathbf{H}_I = \sum_{x_I \in \sup(\mathbf{I})} \mathbf{I}(x_I) * \{\{\{i = x_I.i\} \rightarrow \{\{x_I \rightarrow 1\}\}\}\}$$

$$M_{cp} = \sum_{x_S \in \sup(\mathbf{S})} \sum_{x_R \in \sup(\mathbf{H}_R(\{s=x_S.s\}))} \sum_{x_I \in \sup(\mathbf{H}_I(\{i=x_S.i\}))}$$

$$\mathbf{S}(x_S) * \mathbf{H}_R(\{s = x_S.s\})(x_R) * \mathbf{H}_I(\{i = x_S.i\})(x_I) * x_R.c * x_I.p$$

Transformation Rule: **Aggregate Pushdown**

$$\mathbf{H}_R = \sum_{x_R \in \sup(\mathbf{R})} \mathbf{R}(x_R) * \{\{\{s = x_R.s\} \rightarrow \{\{x_R \rightarrow 1\}\}\}\}$$

$$\mathbf{H}_I = \sum_{x_I \in \sup(\mathbf{I})} \mathbf{I}(x_I) * \{\{\{i = x_I.i\} \rightarrow \{\{x_I \rightarrow 1\}\}\}\}$$

$$M_{cp} = \sum_{x_S \in \sup(\mathbf{S})} \sum_{x_R \in \sup(\mathbf{H}_R(\{s = x_S.s\}))} \sum_{x_I \in \sup(\mathbf{H}_I(\{i = x_S.i\}))}$$

$$\mathbf{S}(x_S) * \mathbf{H}_R(\{s = x_S.s\})(x_R) * \mathbf{H}_I(\{i = x_S.i\})(x_I) * x_R.c * x_I.p$$

Push aggregate $\sum_{x_R} x_R.c$ into $\mathbf{H}_R$

Transformation Rule: **Aggregate Pushdown**

$$\mathbf{H}_R = \sum_{x_R \in \text{sup}(\mathbf{R})} \mathbf{R}(x_R) * \{\{\{s = x_R.s\} \rightarrow x_R.c\}\}$$

$$\mathbf{H}_I = \sum_{x_I \in \text{sup}(\mathbf{I})} \mathbf{I}(x_I) * \{\{\{i = x_I.i\} \rightarrow \{\{x_I \rightarrow 1\}\}\}\}$$

$$M_{cp} = \sum_{x_S \in \text{sup}(\mathbf{S})} \sum_{x_R \in \text{sup}(\mathbf{H}_I(\{s=x_S.s\}))} \sum_{x_I \in \text{sup}(\mathbf{H}_I(\{i=x_S.i\}))}$$

$$\mathbf{S}(x_S) * \mathbf{H}_R(\{s = x_S.s\})(x_R) * \mathbf{H}_I(\{i = x_S.i\})(x_I) * x_I.p$$

Transformation Rule: **Aggregate Pushdown**

$$\mathbf{H}_R = \sum_{x_R \in \text{sup}(\mathbf{R})} \mathbf{R}(x_R) * \{\{\{s = x_R.s\} \to x_R.c\}\}$$

$$\mathbf{H}_I = \sum_{x_I \in \text{sup}(\mathbf{I})} \mathbf{I}(x_I) * \{\{\{i = x_I.i\} \to \{\{x_I \to 1\}\}\}\}$$

$$M_{cp} = \sum_{x_S \in \text{sup}(\mathbf{S})} \sum_{x_I \in \text{sup}(\mathbf{H}_I(\{i = x_S.i\}))}$$

$$\mathbf{S}(x_S) * \mathbf{H}_R(\{s = x_S.s\}) * \mathbf{H}_I(\{i = x_S.i\})(x_I) * x_I.p$$

Transformation Rule: **Aggregate Pushdown**

$$\mathbf{H}_R = \sum_{x_R \in \sup(\mathbf{R})} \mathbf{R}(x_R) * \{\{\{s = x_R.s\} \rightarrow x_R.c\}\}$$

$$\mathbf{H}_I = \sum_{x_I \in \sup(\mathbf{I})} \mathbf{I}(x_I) * \{\{\{i = x_I.i\} \rightarrow \{\{x_I \rightarrow 1\}\}\}\}$$

$$M_{cp} = \sum_{x_S \in \sup(\mathbf{S})} \sum_{x_I \in \sup(\mathbf{H}_I(\{i=x_S.i\}))}$$

$$\mathbf{S}(x_S) * \mathbf{H}_R(\{s = x_S.s\}) * \mathbf{H}_I(\{i = x_S.i\})(x_I) * x_I.p$$

Push aggregate $\sum\limits_{x_I} x_I.p$ into $\mathbf{H}_I$

Transformation Rule: **Aggregate Pushdown**

$$\mathbf{H}_R = \sum_{x_R \in \mathrm{sup}(\mathbf{R})} \mathbf{R}(x_R) * \{\{\{s = x_R.s\} \rightarrow x_R.c\}\}$$

$$\mathbf{H}_I = \sum_{x_I \in \mathrm{sup}(\mathbf{I})} \mathbf{I}(x_I) * \{\{\{i = x_I.i\} \rightarrow x_I.p\}\}$$

$$M_{cp} = \sum_{x_S \in \mathrm{sup}(\mathbf{S})} \sum_{x_I \in \mathrm{sup}(\mathbf{H}_I(\{i=x_S.i\}))}$$

$$\mathbf{S}(x_S) * \mathbf{H}_R(\{s = x_S.s\}) * \mathbf{H}_I(\{i = x_S.i\})$$

Transformation Rule: **Aggregate Pushdown**

$$\mathbf{H}_R = \sum_{x_R \in \sup(\mathbf{R})} \mathbf{R}(x_R) * \{\{\{s = x_R.s\} \rightarrow x_R.c\}\}$$

$$\mathbf{H}_I = \sum_{x_I \in \sup(\mathbf{I})} \mathbf{I}(x_I) * \{\{\{i = x_I.i\} \rightarrow x_I.p\}\}$$

$$M_{cp} = \sum_{x_S \in \sup(\mathbf{S})} \mathbf{S}(x_S) * \mathbf{H}_R(\{s = x_S.s\}) * \mathbf{H}_I(\{i = x_S.i\})$$

Transformation Rule: **Aggregate Pushdown**

$$\mathbf{H}_R = \sum_{x_R \in \sup(\mathbf{R})} \mathbf{R}(x_R) * \{\{\{s = x_R.s\} \rightarrow x_R.c\}\}$$

$$\mathbf{H}_I = \sum_{x_I \in \sup(\mathbf{I})} \mathbf{I}(x_I) * \{\{\{i = x_I.i\} \rightarrow x_I.p\}\}$$

$$M_{cp} = \sum_{x_S \in \sup(\mathbf{S})} \mathbf{S}(x_S) * \mathbf{H}_R(\{s = x_S.s\}) * \mathbf{H}_I(\{i = x_S.i\})$$

---

$$\mathbf{H}'_R = \sum_{x_R \in \sup(\mathbf{R})} \mathbf{R}(x_R) * \{\{\{s = x_R.s\} \rightarrow x_R.c * x_R.c\}\}$$

$$\mathbf{H}'_I = \sum_{x_I \in \sup(\mathbf{I})} \mathbf{I}(x_I) * \{\{\{i = x_I.i\} \rightarrow 1\}\}$$

$$M_{cc} = \sum_{x_S \in \sup(\mathbf{S})} \mathbf{S}(x_S) * \mathbf{H'}_R(\{s = x_S.s\}) * \mathbf{H'}_I(\{i = x_S.i\})$$

Similarly for $M_{cc}$

Transformation Rule: **Aggregate Fusion**

$$\mathbf{H}_R = \sum_{x_R \in \sup(\mathbf{R})} \mathbf{R}(x_R) * \{\{s = x_R.s\} \rightarrow x_R.c\}\} \longleftarrow$$

$$\mathbf{H}'_R = \sum_{x_R \in \sup(\mathbf{R})} \mathbf{R}(x_R) * \{\{s = x_R.s\} \rightarrow x_R.c * x_R.c\}\} \longleftarrow$$

Fuse $\mathbf{H}_R$ and $\mathbf{H}'_R$

Transformation Rule: **Aggregate Fusion**

$$\mathbf{H}_R'' = \sum_{x_R \in \mathrm{sup}(\mathbf{R})} \mathbf{R}(x_R) * \{\{\{s = x_R.s\} \rightarrow \{v_R = x_R.c,\ v_R' = x_R.c * x_R.c\}\}\}$$

$\mathbf{H}_R''$ computes two aggregates

Transformation Rule: **Aggregate Fusion**

$$\mathbf{H}_R'' = \sum_{x_R \in \sup(\mathbf{R})} \mathbf{R}(x_R) * \{\{\{s = x_R.s\} \rightarrow \{v_R = x_R.c, \ v_R' = x_R.c * x_R.c\}\}\}$$

$$\mathbf{H}_I'' = \sum_{x_I \in \sup(\mathbf{I})} \mathbf{I}(x_I) * \{\{\{i = x_I.i\} \rightarrow \{v_I = x_I.p, \ v_I' = 1\}\}\}$$
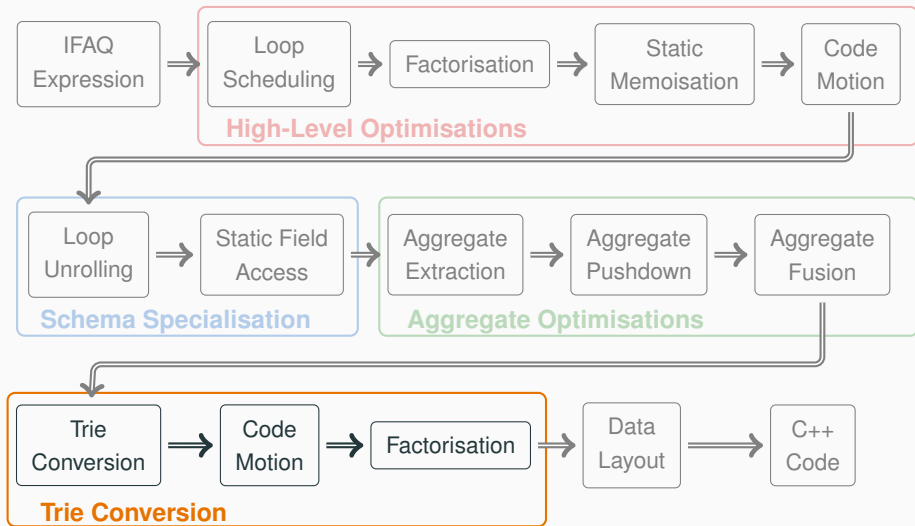
Fuse $\mathbf{H}_I$ and $\mathbf{H}_I'$

Transformation Rule: **Aggregate Fusion**

$$\mathbf{H}_R'' = \sum_{x_R \in \sup(\mathbf{R})} \mathbf{R}(x_R) * \{\{\{s = x_R.s\} \rightarrow \{v_R = x_R.c, \;\; v_R' = x_R.c * x_R.c\}\}\}$$

$$\mathbf{H}_I'' = \sum_{x_I \in \sup(\mathbf{I})} \mathbf{I}(x_I) * \{\{\{i = x_I.i\} \rightarrow \{v_I = x_I.p, \;\; v_I' = 1\}\}\}$$

$$M_{cc,cp} = \sum_{x_S \in \sup(\mathbf{S})} \mathbf{S}(x_S) * \Big($$

$$\texttt{let } w_R = H_R''(\{s = x_S.s\}) \texttt{ in}$$

$$\texttt{let } w_I = H_I''(\{i = x_S.i\}) \texttt{ in}$$

$$\{m_{cp} = w_R.v_R * w_I.v_I, \quad m_{cc} = w_R.v_R' * w_I.v_I'\}\Big)$$

IFAQ Expression → **Loop Scheduling** → Factorisation → **Static Memoisation** → **Code Motion**

**High-Level Optimisations**

Loop Unrolling → Static Field Access

**Schema Specialisation**

Aggregate Extraction → Aggregate Pushdown → Aggregate Fusion

**Aggregate Optimisations**

Trie Conversion → Code Motion → Factorisation → Data Layout → C++ Code

**Trie Conversion**

Transformation Rule: **Trie Conversion**

$$H_R'' = \sum_{x_r \in \sup(\mathbf{R})} \mathbf{R}(x_r) * \{\{\{s = x_r.s\} \rightarrow \{v_R = x_r.c, v_R' = x_r.c * x_r.c\}\}\}$$

$$H_I'' = \sum_{x_i \in \sup(\mathbf{I})} \mathbf{I}(x_i) * \{\{\{i = x_i.i\} \rightarrow \{v_I = x_i.p, v_R' = 1\}\}\}$$

$$M_{cc,cp} = \sum_{x_S \in \sup(\mathbf{S})} \mathbf{S}(x_S) * \Big($$

$$\texttt{let } w_R = H_R''(\{s = x_S.s\}) \texttt{ in}$$

$$\texttt{let } w_I = H_I''(\{i = x_S.i\}) \texttt{ in}$$

$$\{m_{cp} = w_R.v_R * w_I.v_I, \quad m_{cc} = w_R.v_R' * w_I.v_I'\}\Big)$$

Turn relations into tries (i.e., nested dictionaries)

Transformation Rule: **Trie Conversion**

$$M_{cc,cp} = \sum_{x_S \in \text{sup}(\mathbf{S})} \mathbf{S}(x_S) * \Big($$

$$\texttt{let } w_R = \mathbf{H}_R''(\{s = x_S.s\}) \texttt{ in}$$

$$\texttt{let } w_I = \mathbf{H}_I''(\{i = x_S.i\}) \texttt{ in}$$

$$\{m_{cp} = w_R.v_R * w_I.v_I, \quad m_{cc} = w_R.v_R' * w_I.v_I'\}\Big)$$

Turn relation S into trie S'

Transformation Rule: **Trie Conversion**

One loop for each join variable

$$M_{cc,cp} = \sum_{x_S \in \mathrm{sup}(\mathbf{S'})}$$

$$\sum_{x_i \in \mathrm{sup}(\mathbf{S'}(x_s))} S'(x_s)(x_i) * \Big($$

$$\mathtt{let}\ \ w_R = \mathbf{H}''_R(\{s = x_S.s\})\ \mathtt{in}$$

$$\mathtt{let}\ \ w_I = \mathbf{H}''_I(\{i = x_S.i\})\ \mathtt{in}$$

$$\{m_{cp} = w_R.v_R * w_I.v_I, \quad m_{cc} = w_R.v'_R * w_I.v'_I\}\Big)$$

Transformation Rule: **Code Motion**

$$M_{cc,cp} = \sum_{x_S \in \sup(\mathbf{S'})}$$

$$\sum_{x_i \in \sup(\mathbf{S'}(x_s))} S'(x_s)(x_i) * \Big($$

$$\texttt{let } w_R = \mathbf{H}''_R(\{s = x_S.s\}) \texttt{ in}$$

$$\texttt{let } w_I = \mathbf{H}''_I(\{i = x_S.i\}) \texttt{ in}$$

$$\{m_{cp} = w_R.v_R * w_I.v_I, \quad m_{cc} = w_R.v'_R * w_I.v'_I\}\Big)$$

> Move up the look-up into $\mathbf{H}_R$

Transformation Rule: **Code Motion**

$$M_{cc,cp} = \sum_{x_S \in \sup(\mathbf{S'})}$$

> Move up the look-up into $\mathbf{H}_R$

$$\texttt{let } w_R = \mathbf{H}_R''(\{s = x_S.s\}) \texttt{ in}$$

$$\sum_{x_i \in \sup(\mathbf{S'}(x_S))} S'(x_S)(x_i) * \Big($$

$$\texttt{let } w_I = \mathbf{H}_I''(\{i = x_S.i\}) \texttt{ in}$$

$$\{m_{cp} = w_R.v_R * w_I.v_I, \quad m_{cc} = w_R.v_R' * w_I.v_I'\}\Big)$$

Transformation Rule: **Factorisation**

$$M_{cc,cp} = \sum_{x_S \in \sup(\mathbf{S'})}$$

$$\texttt{let } w_R = \mathbf{H}_R''(\{s = x_S.s\}) \texttt{ in}$$

$$\sum_{x_i \in \sup(\mathbf{S'}(x_S))} S'(x_s)(x_i) * \Big($$
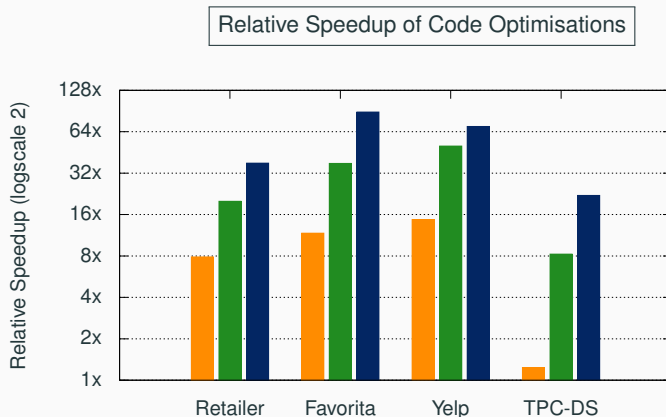
$$\texttt{let } w_I = \mathbf{H}_I''(\{i = x_S.i\}) \texttt{ in}$$

$$\{m_{cp} = w_R.v_R * w_I.v_I, \quad m_{cc} = w_R.v_R' * w_I.v_I'\} \Big)$$

Less arithmetic operations

Transformation Rule: **Factorisation**

Less arithmetic operations

$$M_{cc,cp} = \sum_{x_S \in \text{sup}(\mathbf{S'})}$$

$$\texttt{let } w_R = \mathbf{H}_R''(\{s = x_S.s\}) \texttt{ in}$$

$$\{m_{cp} = w_R.v_R, \quad m_{cc} = w_R.v_R'\} *$$

$$\sum_{x_i \in \text{sup}(\mathbf{S'}(x_s))} S'(x_s)(x_i) * \Big($$

$$\texttt{let } w_I = \mathbf{H}_I''(\{i = x_S.i\}) \texttt{ in}$$

$$\{m_{cp} = w_I.v_I, \quad m_{cc} = w_I.v_I'\}\Big)$$

Relative Speedup of Code Optimisations

Added optimisations for covariance matrix computation:

specialisation $\rightarrow$ + sharing $\rightarrow$ + parallelisation

AWS d2.xlarge (4 vCPUs, 32GB)

Three-step recipe for efficient machine learning over databases:

1. Turn the learning problem into a database problem

2. Exploit the problem structure to lower the complexity

3. Specialise and optimise the code to lower the constant factors

Q.E.D.