

On Extending PostgreSQL with the Skyline Operator

Hannes Eder

Technische Universität Wien
Institut für Informationssysteme
Database and Artificial Intelligence Group
Betreuer: Univ.-Prof. Mag.rer.nat. Dr.techn. Reinhard Pichler
Univ.Ass. Dr.rer.nat. Fang Wei

Motivation and Goal

In the realm of multi-criteria optimization designing an appropriate objective function is a challenging task from a user's perspective. The *skyline operator* filters out the *interesting* tuples from a potentially large dataset. A tuple belongs to the *skyline* if it is *not dominated* by any other tuple, i.e. there is no tuple which is at least as good as in all and better in at least one criteria. No matter how we weight our preferences along the attributes, only those tuples which score best under a *monotone scoring function* are part of the skyline. In other words, the skyline does not contain tuples which are nobody's favorite. The notion of skyline is also called *Pareto optimal set* and its computation *maximum vector problem*.

In this thesis we aim at:

- extending PostgreSQL with the skyline operator,
- evaluating skyline algorithms in the RDBMS context, and
- building a skyline query optimizer to automatically generate a good query plan w.r.t. I/O, time, and memory consumption.

Formal Problem Definition

In a skyline query with the following skyline clause

`SKYLINE OF a1 MIN, ..., ak MIN, ak+1 MAX, ..., al MAX, al+1 DIFF, ..., am DIFF`

a tuple r

$$r = (r_1, \dots, r_k, \underbrace{r_{k+1}, \dots, r_l}_{\text{MIN}}, \underbrace{r_{l+1}, \dots, r_m}_{\text{DIFF}}, \underbrace{r_{m+1}, \dots, r_n}_{\text{extra}})$$

dominates a tuples s

$$s = (\underbrace{s_1, \dots, s_k}_{\text{MIN}}, \underbrace{s_{k+1}, \dots, s_l}_{\text{MAX}}, \underbrace{s_{l+1}, \dots, s_m}_{\text{DIFF}}, \underbrace{s_{m+1}, \dots, s_n}_{\text{extra}})$$

iff

$$C(r, s) = \left(\bigwedge_{1 \leq i \leq k} r_i \leq s_i \right) \wedge \left(\bigwedge_{k+1 \leq i \leq l} r_i \geq s_i \right) \wedge \left(\bigwedge_{l+1 \leq i \leq m} r_i = s_i \right) \wedge \left(\left(\bigvee_{1 \leq i \leq k} r_i < s_i \right) \vee \left(\bigvee_{k+1 \leq i \leq l} r_i > s_i \right) \right).$$

We define the skyline of a relation R as

$$\text{skyline}_>(R) := \{r \in R \mid \nexists s \in R: s > r\}$$

with

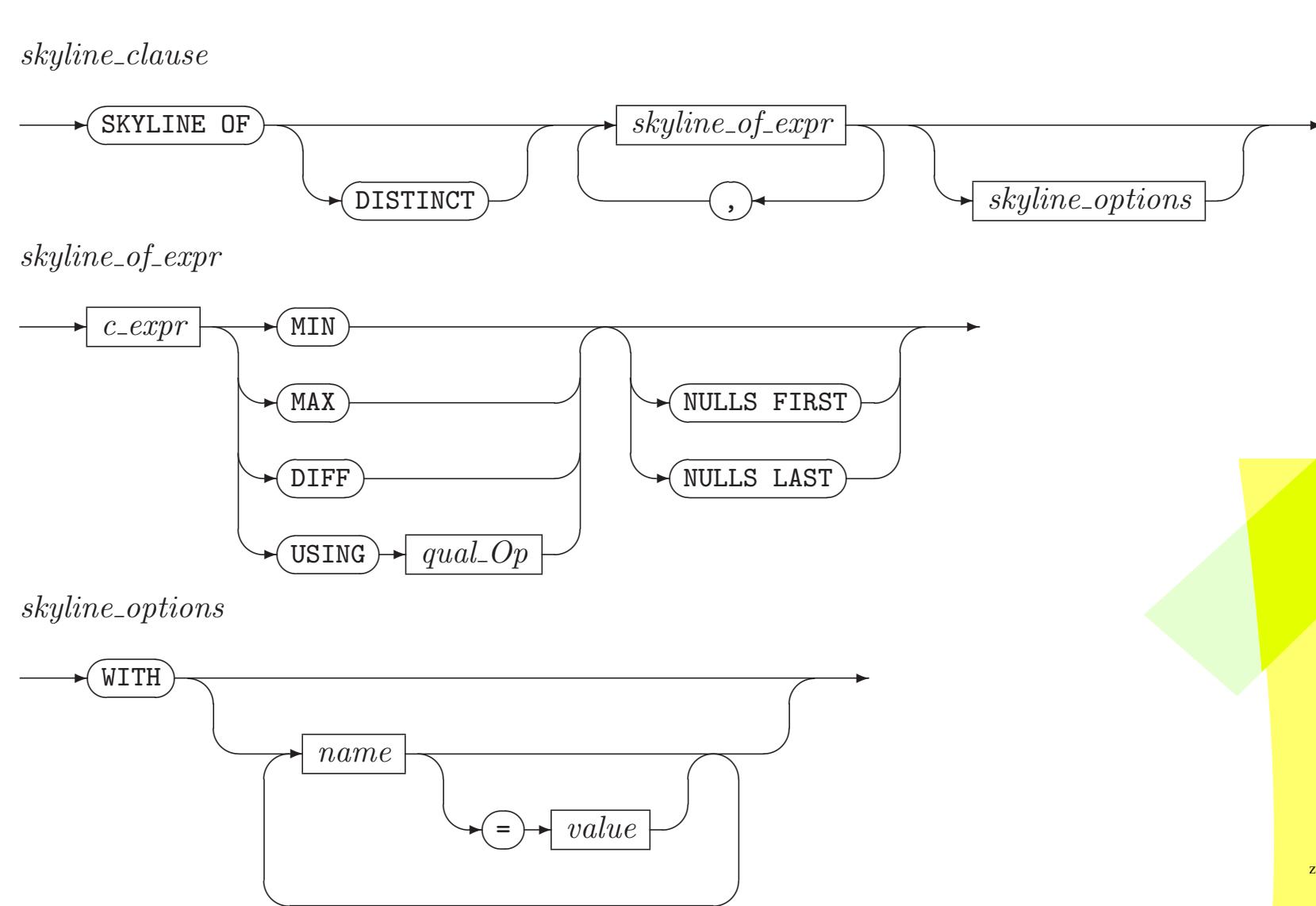
$$r > s \quad \text{iff} \quad C(r, s).$$

Query Syntax

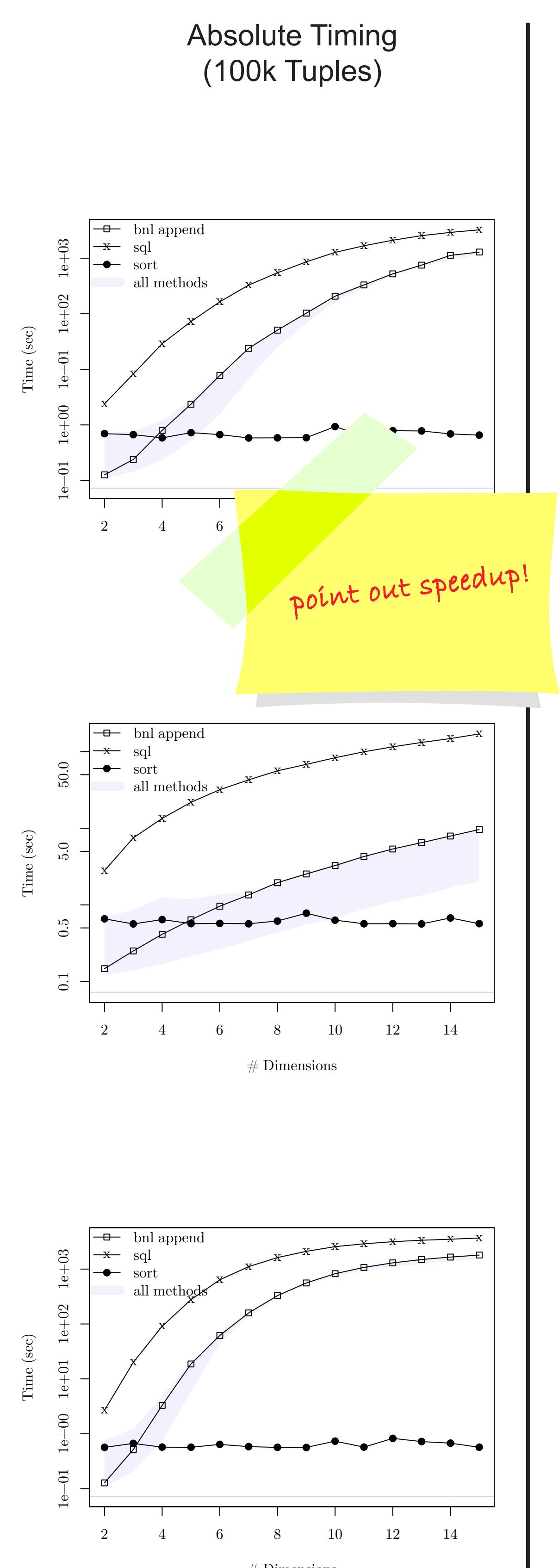
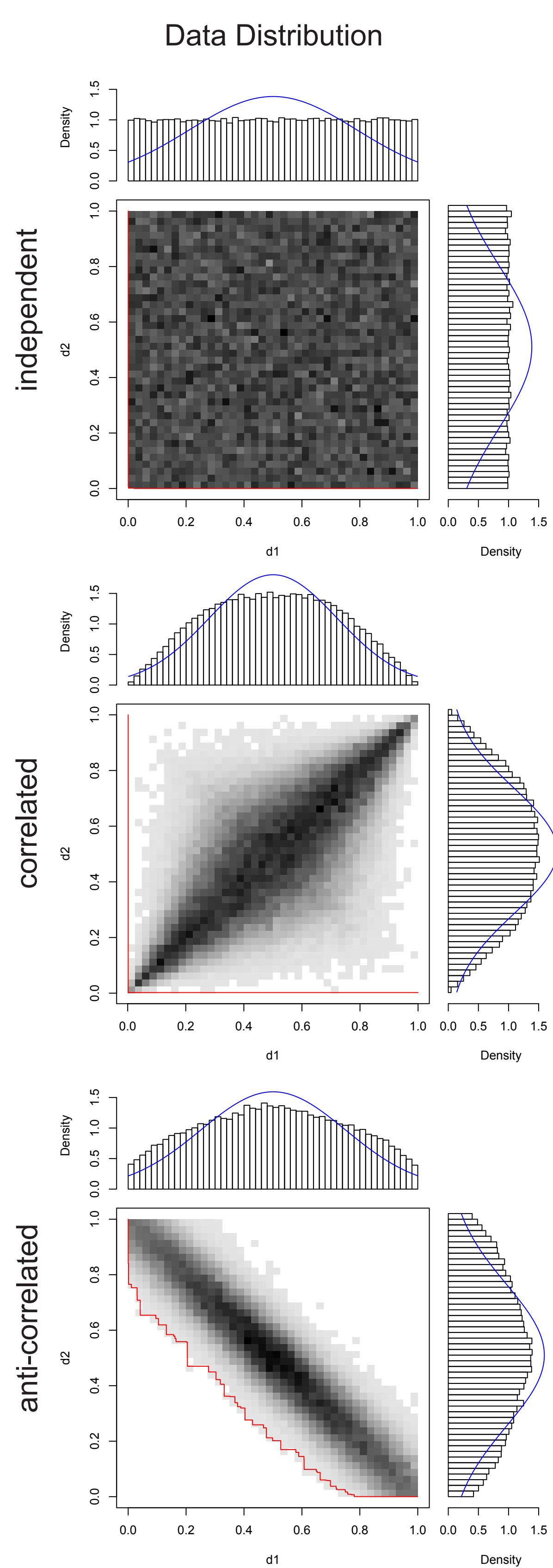
We extend the syntax from Börzsönyi et. al. [2001]

```
SELECT ... FROM ... WHERE ...
GROUP BY ... HAVING ...
SKYLINE OF [DISTINCT] a, [MIN|MAX|DIFF], ... ,
          am [MIN|MAX|DIFF]
ORDER BY ...
```

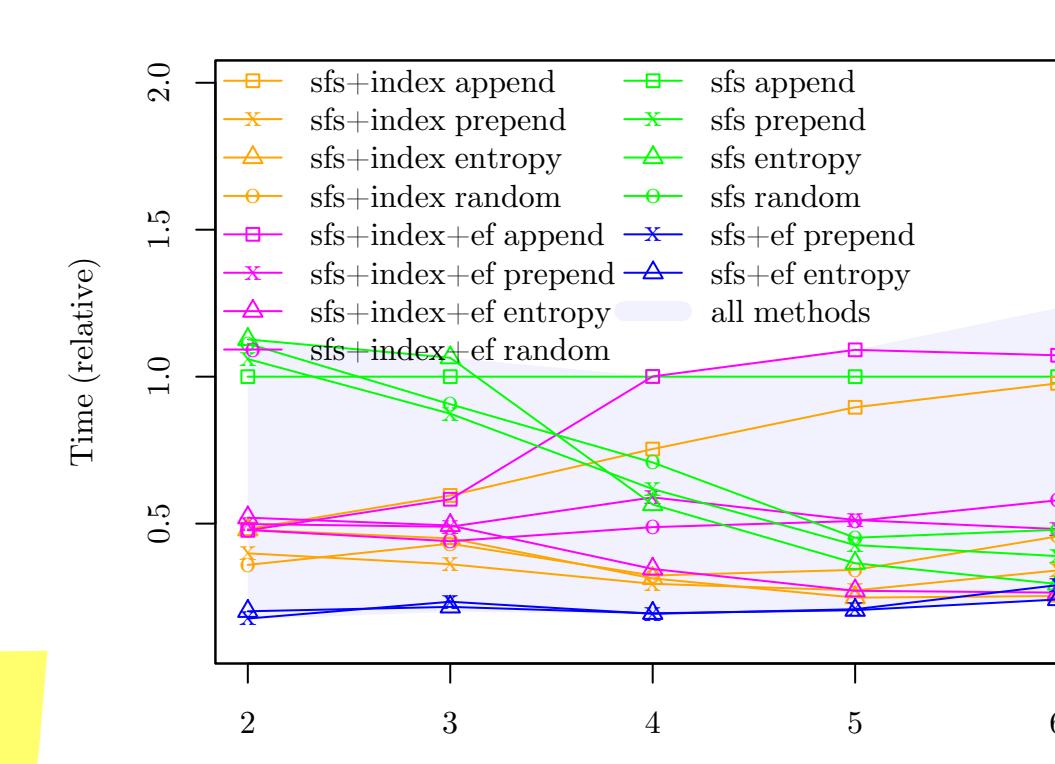
in the following way:



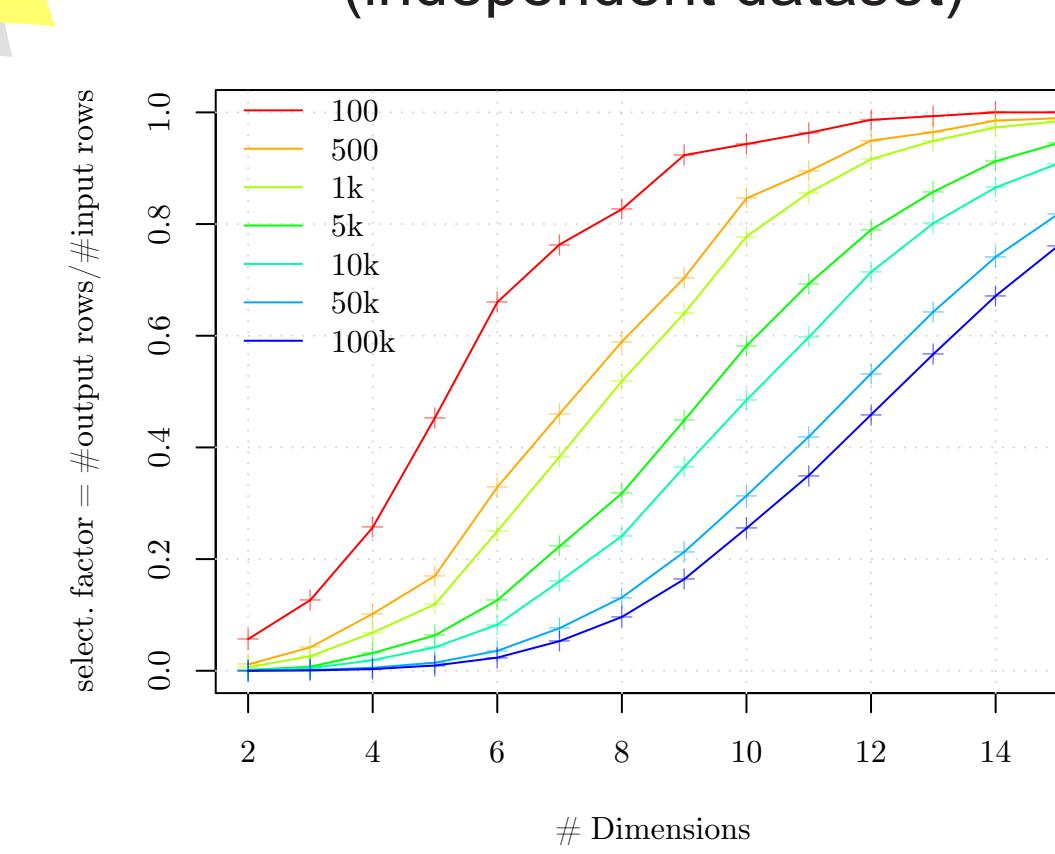
Experimental Results



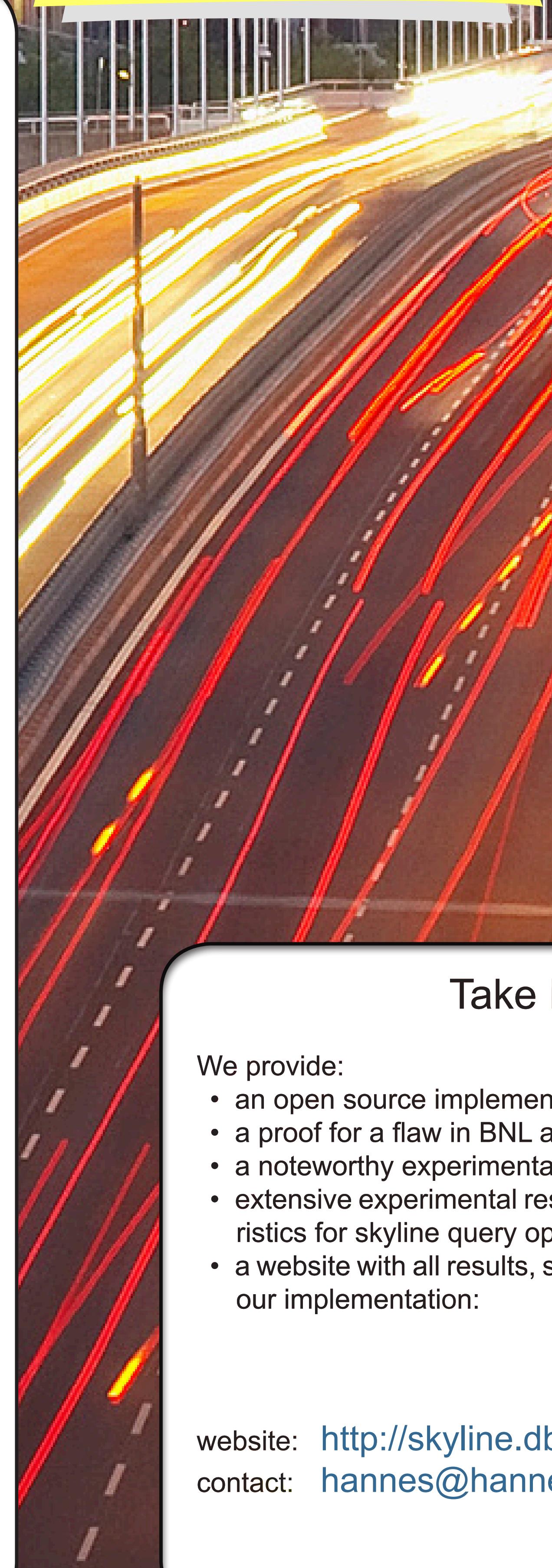
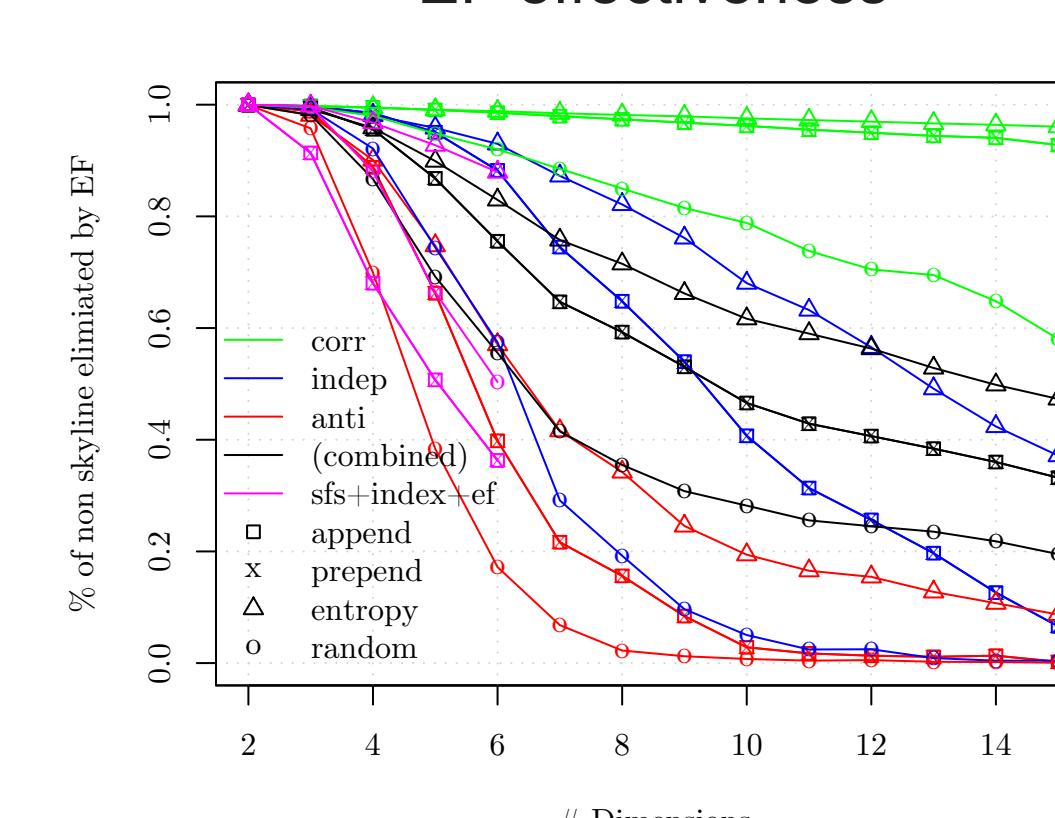
Relative Timing w/ index (100k tuples)



Selectivity Factor (independent dataset)

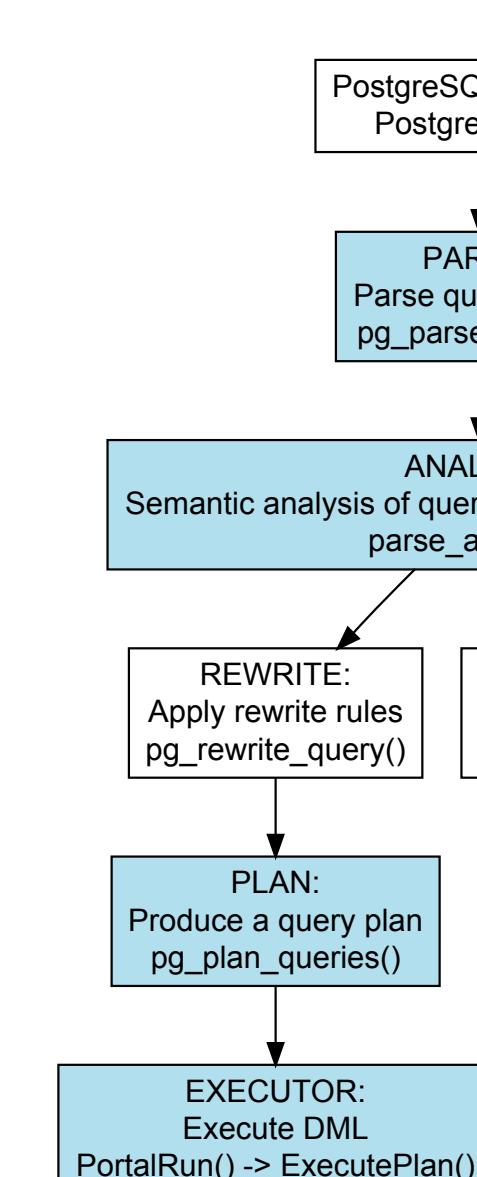


EF effectiveness



Implementation

In this architecture diagram of PostgreSQL the light blue background indicates the components we have modified to support skyline queries.



Take Home Message

We provide:

- an open source implementation of BNL and SFS in PostgreSQL
- a proof for a flaw in BNL and a fixed version
- a noteworthy experimental environment
- extensive experimental results which are beneficial for developing heuristics for skyline query optimization
- a website with all results, source code, and a web-interface to test drive our implementation:

website: <http://skyline.dbai.tuwien.ac.at/>
contact: hannes@hanneseder.net

