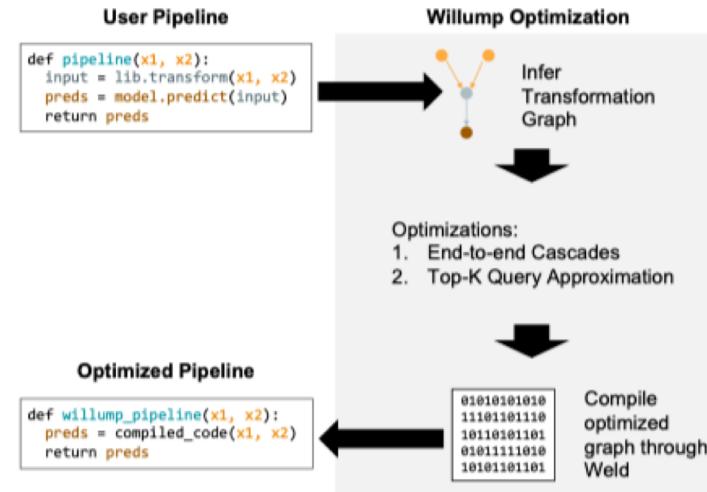


Willump

A Statistically-Aware End-to-end Optimizer for Machine Learning Inference

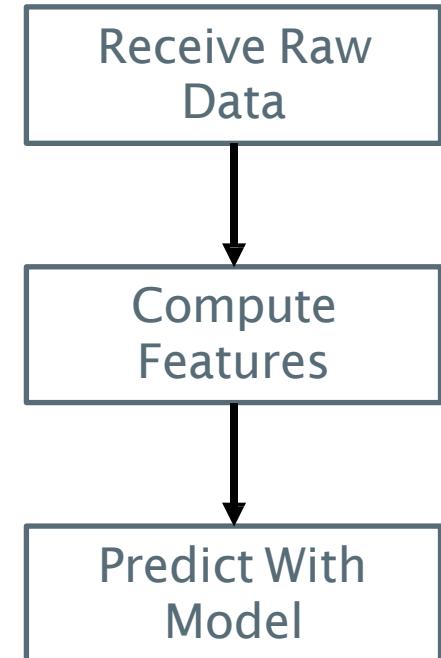


Problem: ML Inference

- Often performance-critical.
- Recent focus on tools for ML prediction serving.

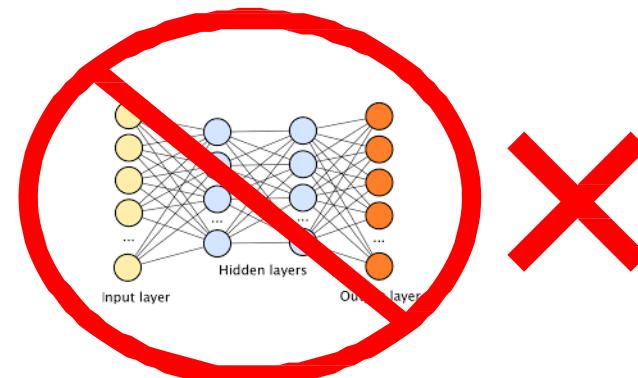
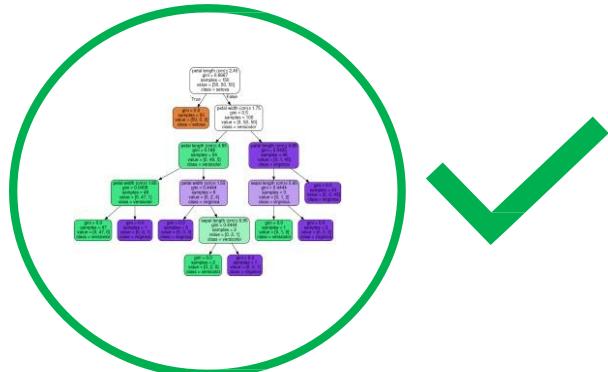
A Common Bottleneck: Feature Computation

- Many applications bottlenecked by feature computation.
- Pipeline of transformations computes numerical *features* from data for model.



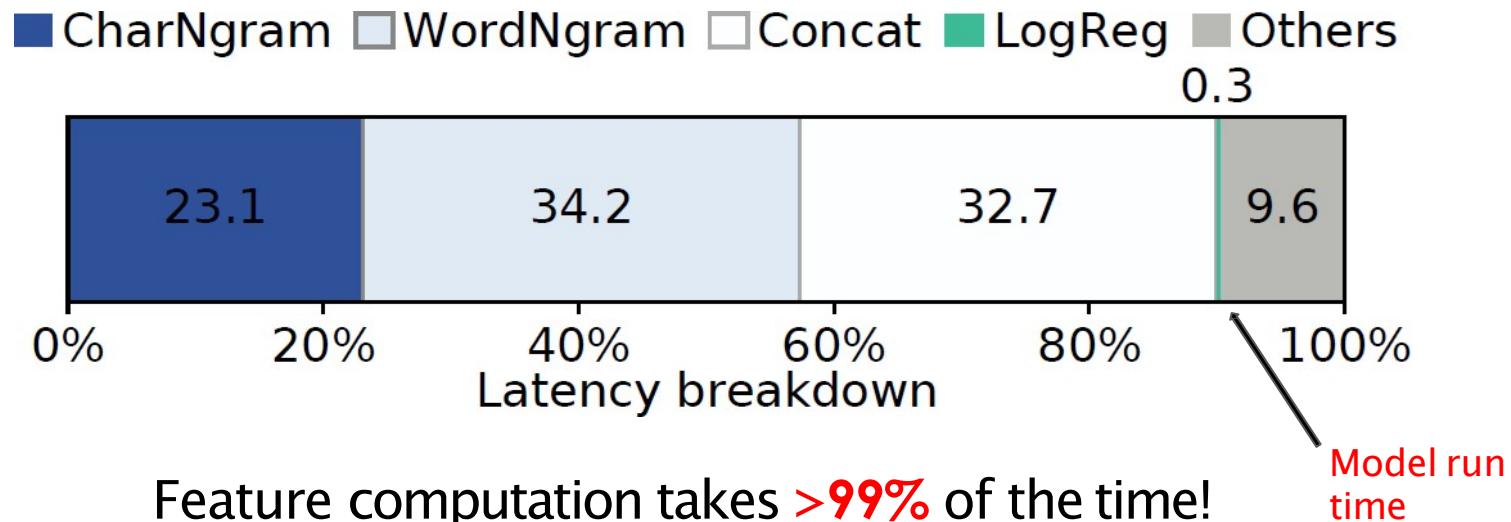
A Common Bottleneck: Feature Computation

- Feature computation is bottleneck when models are inexpensive—boosted trees, not DNNs.
- Common on tabular/structured data!



A Common Bottleneck: Feature Computation

Production Microsoft sentiment analysis pipeline



Current State-of-the-art

- Apply traditional serving optimizations, e.g. caching (Clipper), compiler optimizations (Pretzel).
- Neglect unique **statistical properties** of ML apps.



Statistical Properties of ML

Amenability to approximation

Statistical Properties of ML

Amenability to approximation



Easy input:
Definitely not
a dog.



Hard input:
Maybe a
dog?

Statistical Properties of ML

Amenability to approximation



Easy input:
Definitely not
a dog.



Hard input:
Maybe a
dog?

Existing Systems: Use Expensive Model for Both

Statistical Properties of ML

Amenability to approximation



Easy input:
Definitely not
a dog.



Hard input:
Maybe a
dog?

Statistically-Aware Systems: Use cheap model on bucket,
expensive model on cat.

Statistical Properties of ML

- Model is often part of a bigger app (e.g. top-K query)

Statistical Properties of ML

- Model is often part of a bigger app (e.g. top-K query)

Artist	Score	Rank
Beatles	9.7	1
Bruce Springsteen	9.5	2
...
Justin Bieber	5.6	999
Nickelback	4.1	1000

Problem:
Return top
10 artists.

Statistical Properties of ML

- Model is often part of a bigger app (e.g. top-K query)

Existing Systems

Artist	Score	Rank
Beatles	9.7	1
Bruce Springsteen	9.5	2
...
Justin Bieber	5.6	999
Nickelback	4.1	1000

Use
expensive
model for
everything!

Statistical Properties of ML

- Model is often part of a bigger app (e.g. top-K query)

Statistically-aware Systems

Artist	Score	Rank
Beatles	9.7	1
Bruce Springsteen	9.5	2
...
Justin Bieber	5.6	999
Nickelback	4.1	1000

High-value:
Rank precisely,
return.

Low-value:
Approximate,
discard.

Prior Work: Statistically-Aware Optimizations

- Statistically-aware optimizations exist in literature.
- Always application-specific and custom-built.
- Never automatic!

Source: Cheng et al. (DLRS' 16),
Kang et al. (VLDB '17)

ML Inference Dilemma

- ML inference systems:
 - Easy to use.
 - Slow.
- Statistically-aware systems:
 - Fast
 - Require a lot of work to implement.

Can an ML inference system be fast and easy to use?

Willump: Overview

- Statistically-aware optimizer for ML Inference.
- Targets feature computation!
- Automatic model-agnostic statistically-aware opts.
- 10x throughput+latency improvements.

Outline

- **System Overview**
- Optimization 1: End-to-end Cascades
- Optimization 2: Top-K Query Approximation
- Evaluation

Willump: Goals

- Automatically maximize performance of ML inference applications whose performance bottleneck is feature computation

System Overview

Input Pipeline

```
def pipeline(x1, x2):  
    input = lib.transform(x1, x2)  
    preds = model.predict(input)  
    return preds
```

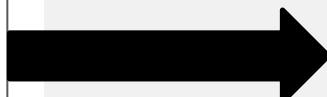
System Overview

Input Pipeline

```
def pipeline(x1, x2):  
    input = lib.transform(x1, x2)  
    preds = model.predict(input)  
    return preds
```

Willump Optimization

Infer Transformation
Graph

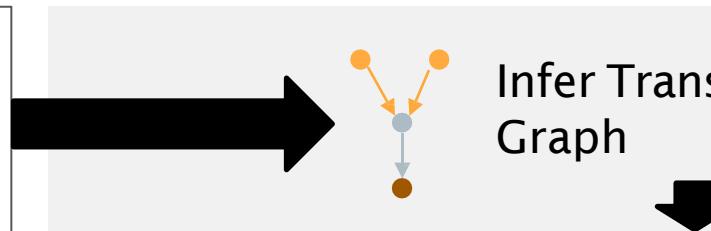


System Overview

Input Pipeline

```
def pipeline(x1, x2):  
    input = lib.transform(x1, x2)  
    preds = model.predict(input)  
    return preds
```

Willump Optimization



Infer Transformation
Graph

Statistically-Aware Optimizations:
1. End-To-End Cascades
2. Top-K Query Approximation

System Overview

Input Pipeline

```
def pipeline(x1, x2):  
    input = lib.transform(x1, x2)  
    preds = model.predict(input)  
    return preds
```

Willump Optimization



Infer Transformation
Graph



Statistically-Aware Optimizations:
1. End-To-End Cascades
2. Top-K Query Approximation



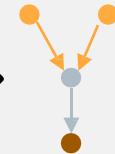
Compiler Optimizations
(Weld—Palkar et al. VLDB '18)

System Overview

Input Pipeline

```
def pipeline(x1, x2):  
    input = lib.transform(x1, x2)  
    preds = model.predict(input)  
    return preds
```

Willump Optimization



Infer Transformation
Graph



Statistically-Aware Optimizations:
1. End-To-End Cascades
2. Top-K Query Approximation

Optimized Pipeline

```
def willump_pipeline(x1, x2):  
    preds = compiled_code(x1, x2)  
    return preds
```

Compiler Optimizations
(Weld—Palkar et al. VLDB '18)



Outline

- System Overview
- **Optimization 1: End-to-end Cascades**
- Optimization 2: Top-K Query Approximation
- Evaluation

Background: Model Cascades

- Classify “easy” inputs with cheap model.
- *Cascade* to expensive model for “hard” inputs.



Easy input:
Definitely not
a dog.



Hard input:
Maybe a
dog?

Background: Model Cascades

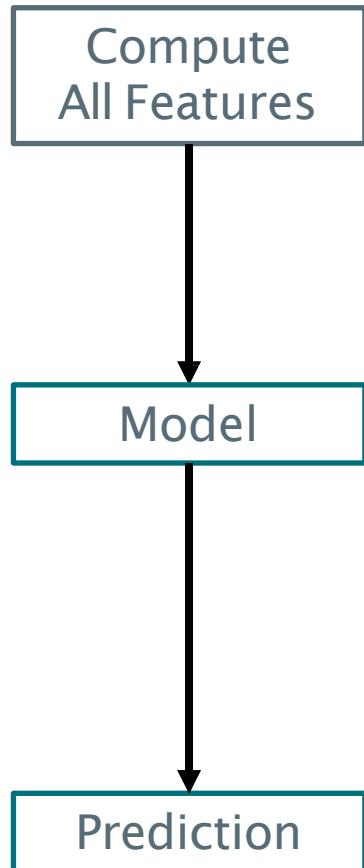
- Used for image classification, object detection.
- Existing systems application-specific and custom-built.

Source: Viola-Jones (CVPR' 01),
Kang et al. (VLDB '17)

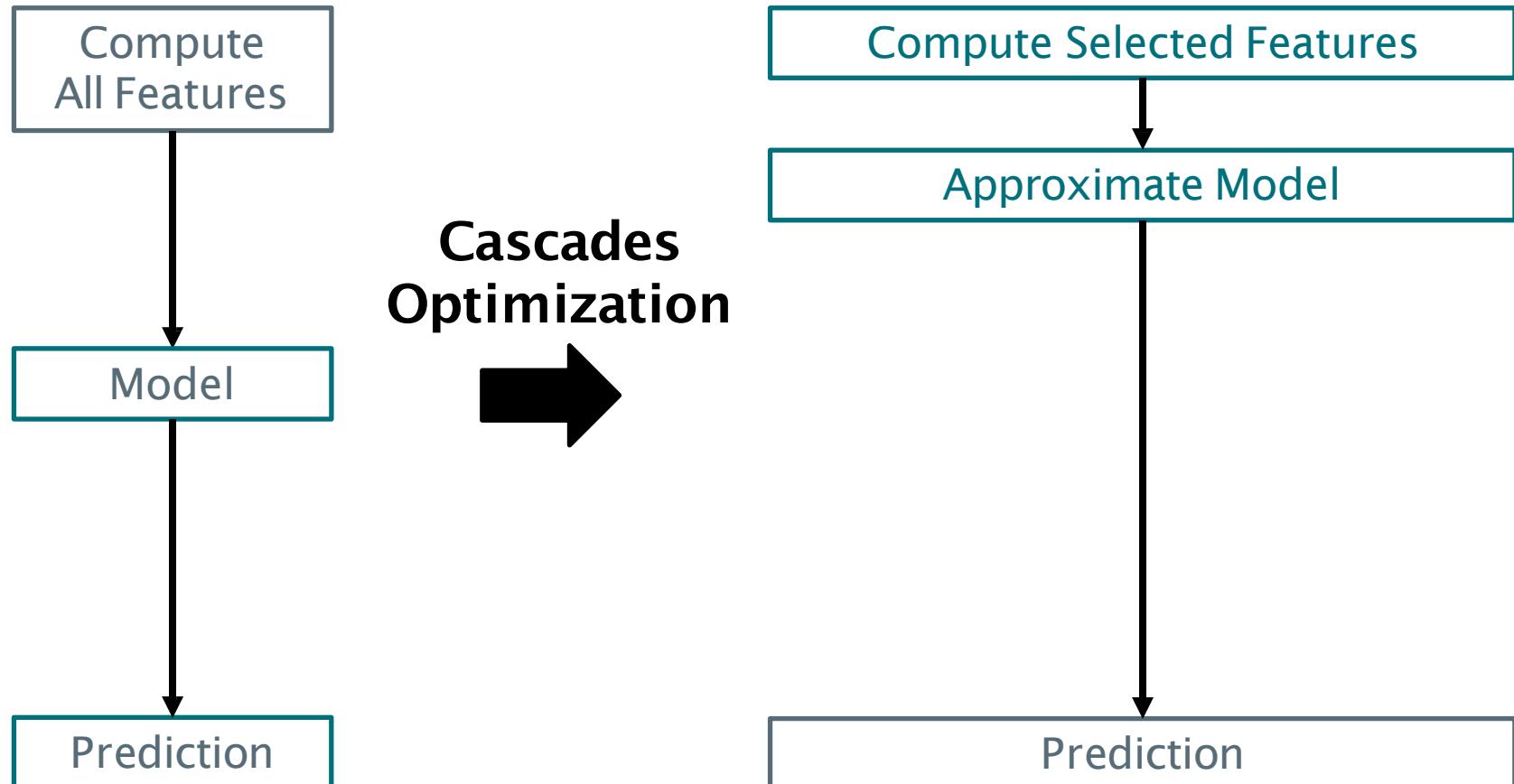
Our Optimization: End-to-end cascades

- Compute only some features for “easy” data inputs; cascade to computing all for “hard” inputs.
- Automatic and model-agnostic, unlike prior work.
 - Estimates for runtime performance & accuracy of a feature set
 - Efficient search process for tuning parameters

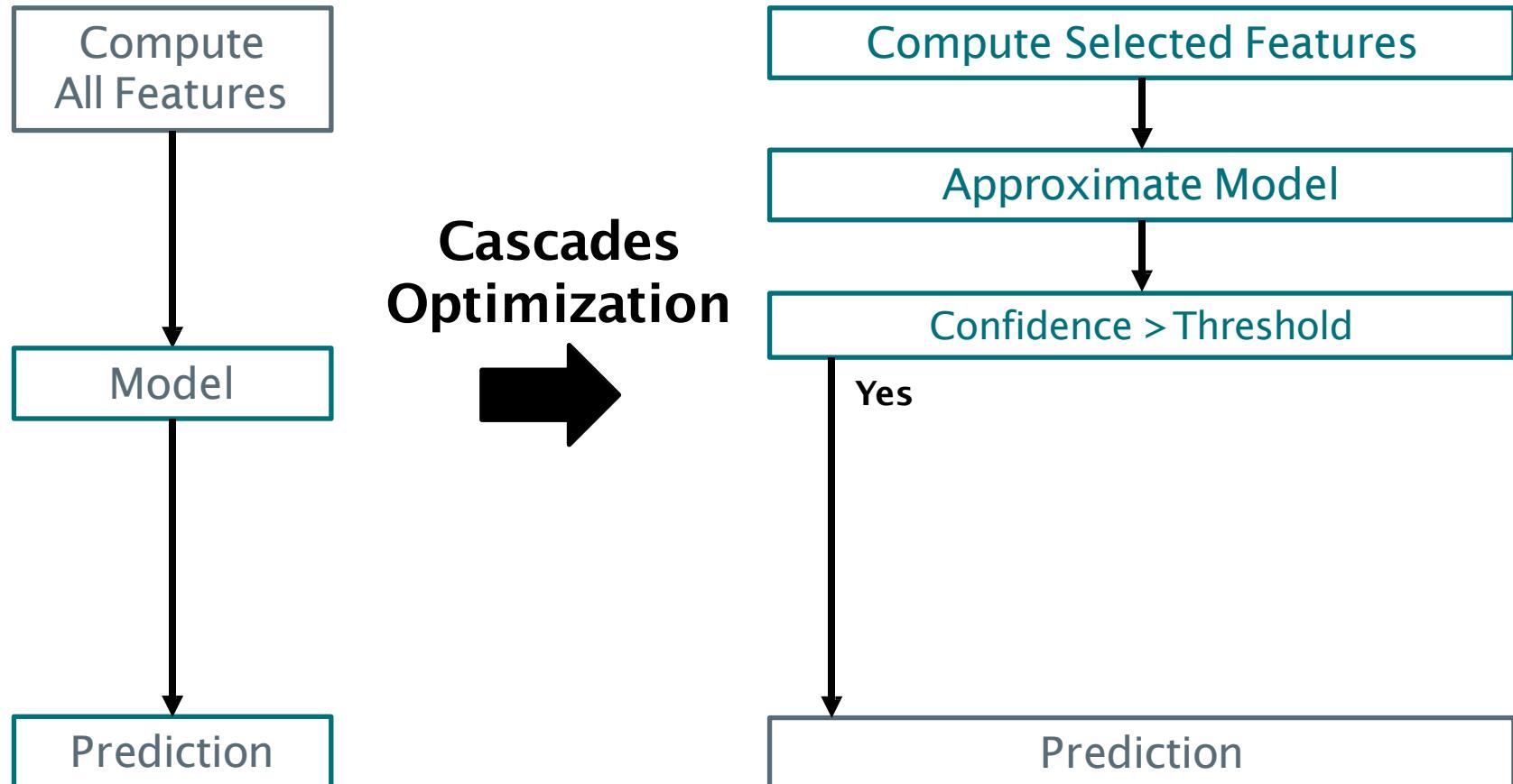
End-to-end Cascades: Original Model



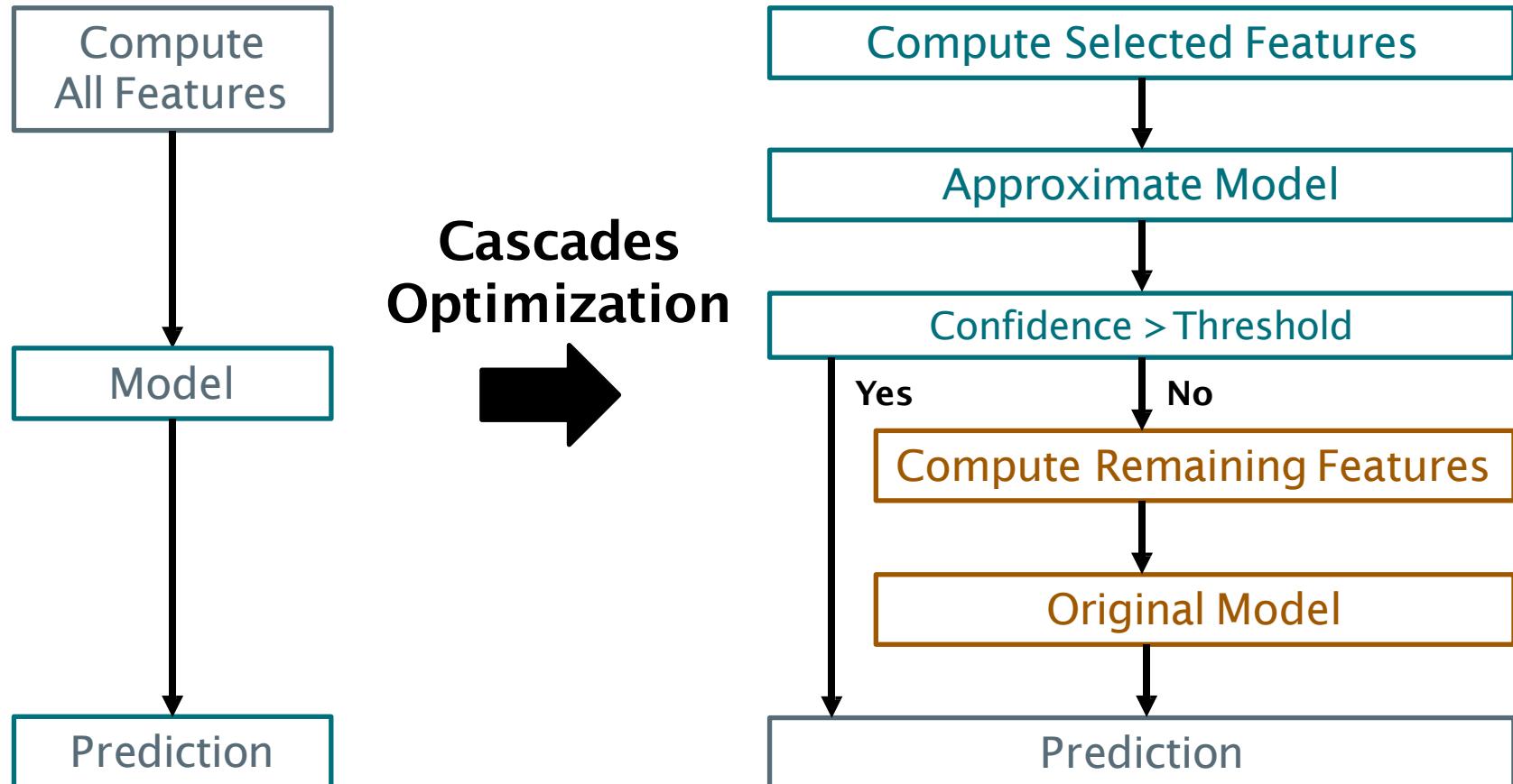
End-to-end Cascades: Approximate Model



End-to-end Cascades: Confidence



End-to-end Cascades: Final Pipeline

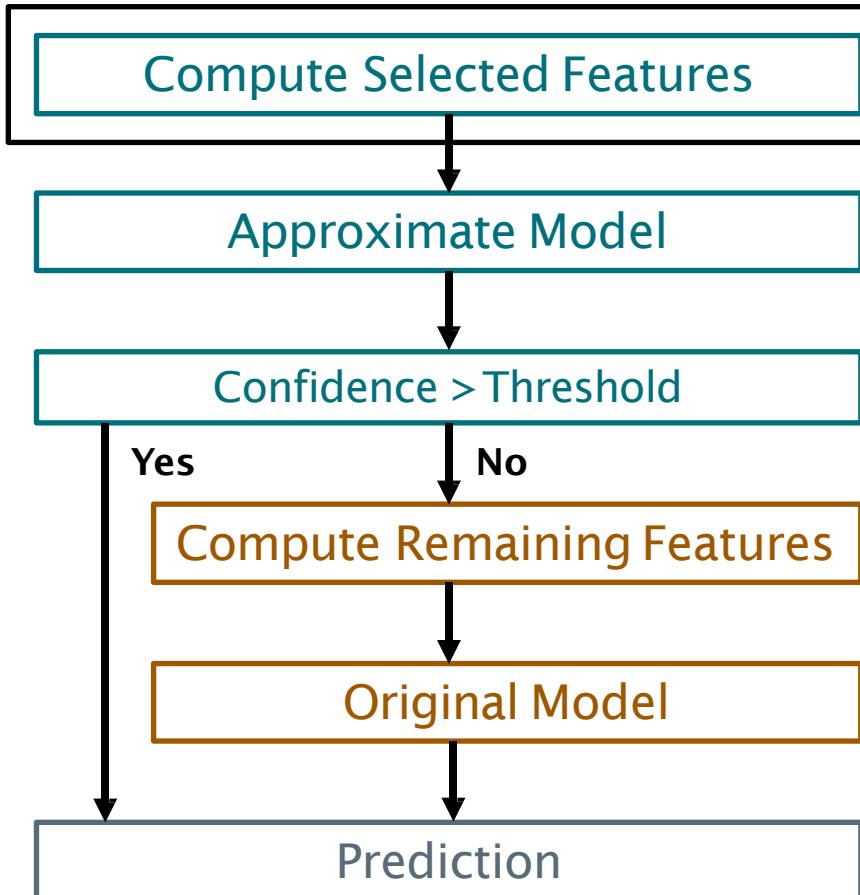


End-to-end Cascades: Constructing Cascades

- Construct cascades during model training.
- Need model training set and an accuracy target.

End-to-end Cascades: Selecting Features

Key question:
Select which
features?

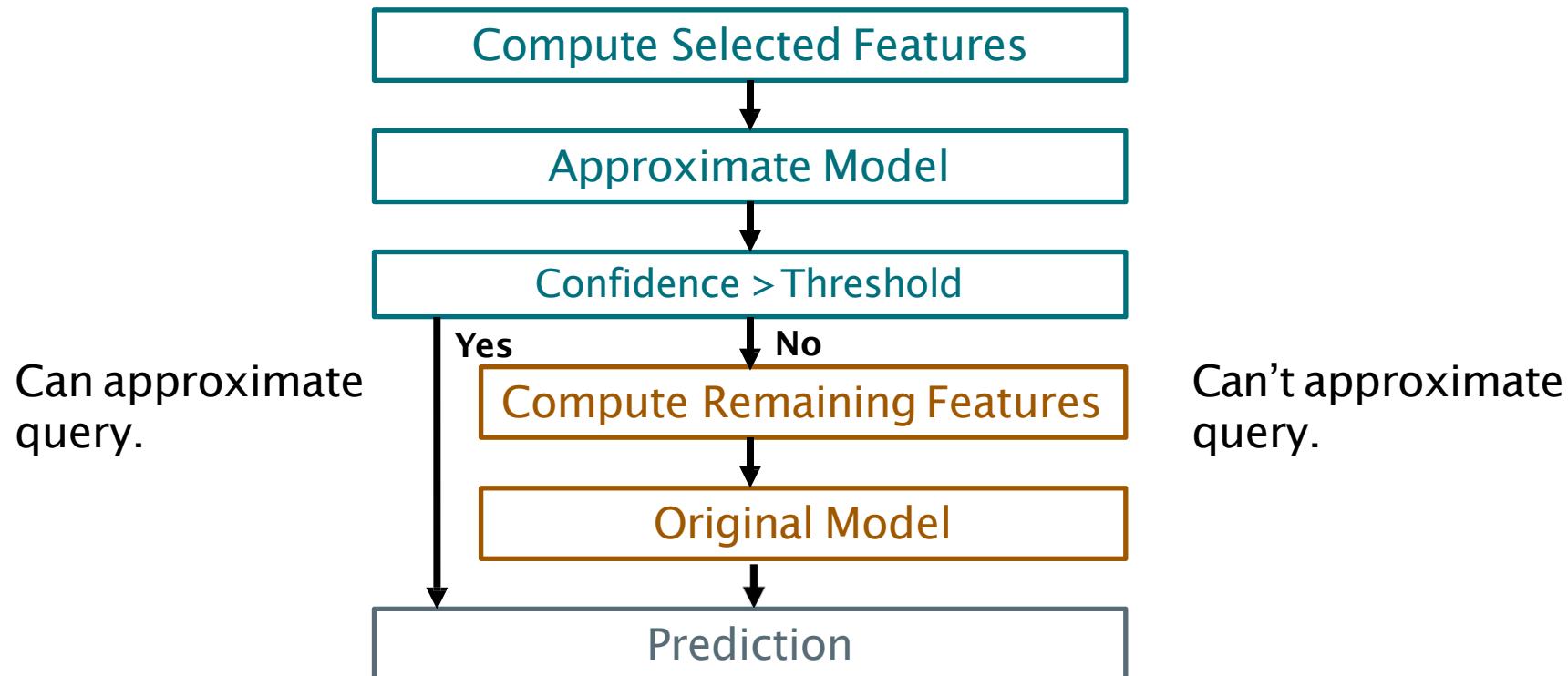


End-to-end Cascades: Selecting Features

- Goal: Select features that minimize expected querytime given accuracy target.

End-to-end Cascades: Selecting Features

Two possibilities for a query: Can approximate or not.

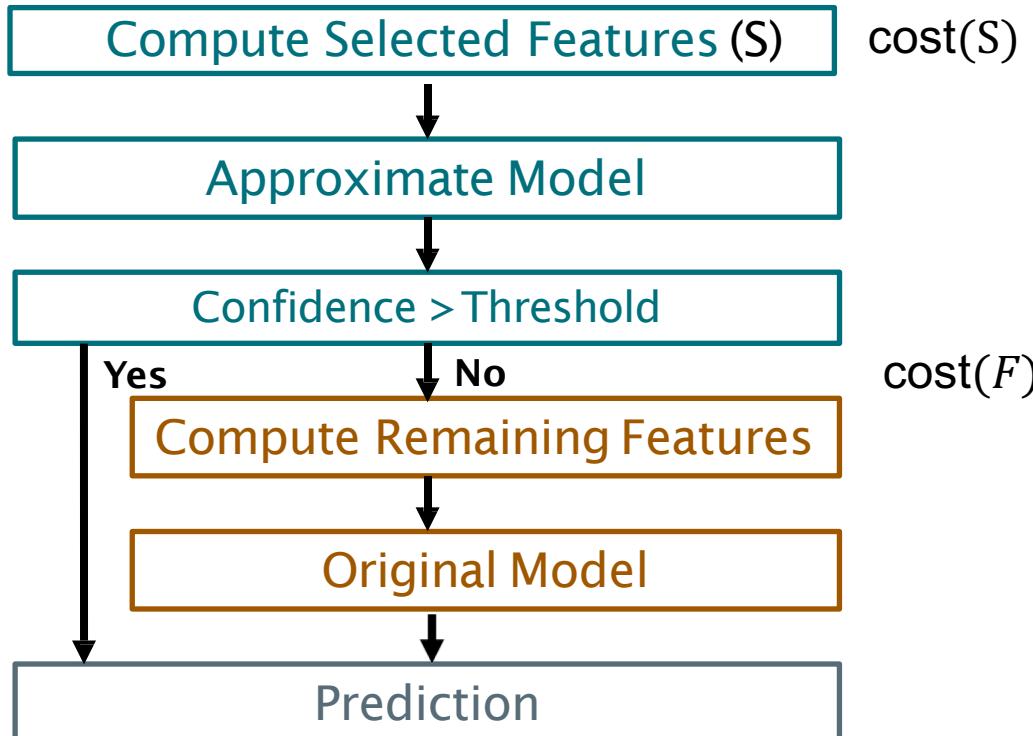


End-to-end Cascades: Selecting Features

$$\min_{S} h_s \cdot \text{cost}(S) + (1-h_s) \cdot \text{cost}(F)$$

h_s

h_s is the % of the features used in the approx model



End-to-end Cascades: Selecting Features

- Goal: Select feature set S that minimizes query time:

$$\min_S h_s \cdot \text{cost}(S) + (1 - h_s) \cdot \text{cost}(F)$$

End-to-end Cascades: Selecting Features

- Goal: Select feature set S that minimizes query time:

$$\min_S h_s \cdot \text{cost}(S) + (1 - h_s) \cdot \text{cost}(F)$$

- Approach:
 - **Choose several potential values of $\text{cost}(S)$.**
 - Find best feature set with each $\text{cost}(S)$.
 - Train model & find cascade threshold for each set.
 - Pick best overall.

End-to-end Cascades: Selecting Features

- Goal: Select feature set S that minimizes query time:

$$\min_S h_s \cdot \text{cost}(S) + (1 - h_s) \cdot \text{cost}(F)$$

- Approach:
 - Choose several potential values of $\text{cost}(S)$.
 - **Find best feature set with each $\text{cost}(S)$.**
 - Train model & find cascade threshold for each set.
 - Pick best overall.

End-to-end Cascades: Selecting Features

- Goal: Select feature set S that minimizes query time:

$$\min_S h_s \cdot \text{cost}(S) + (1 - h_s) \cdot \text{cost}(F)$$

- Approach:
 - Choose several potential values of $\text{cost}(S)$.
 - Find best feature set with each $\text{cost}(S)$.
 - **Train model & find cascade threshold for each set.**
 - Pick best overall.

End-to-end Cascades: Selecting Features

- Goal: Select feature set S that minimizes query time:

$$\min_S h_s \cdot \text{cost}(S) + (1 - h_s) \cdot \text{cost}(F)$$

- Approach:
 - Choose several potential values of $\text{cost}(S)$.
 - Find best feature set with each $\text{cost}(S)$.
 - Train model & find cascade threshold for each set.
 - **Pick best overall.**

End-to-end Cascades: Selecting Features

- Subgoal: Find S minimizing query time if $cost(S) = c_{max}$.

$$\min_{S} h_s \cdot \text{cost}(S) + (1 - h_s) \cdot \text{cost}(F)$$

End-to-end Cascades: Selecting Features

- Subgoal: Find S minimizing query time if $cost(S) = c_{max}$.

$$\min_S h_s \cdot cost(S) + (1 - h_s) \cdot cost(F)$$

- Solution:

- Find S maximizing approximate model accuracy.

End-to-end Cascades: Selecting Features

- Subgoal: Find S minimizing query time if $cost(S) = c_{max}$.

$$\min_S h_s \cdot cost(S) + (1 - h_s) \cdot cost(F)$$

- Solution:

- Find S maximizing approximate model accuracy.
- Problem: Computing accuracy expensive.

End-to-end Cascades: Selecting Features

- Subgoal: Find S minimizing query time if $cost(S) = c_{max}$.

$$\min_S h_s \cdot cost(S) + (1 - h_s) \cdot cost(F)$$

- Solution:

- Find S maximizing approximate model accuracy.
- Problem: Computing accuracy expensive.
- Solution: Estimate accuracy via **permutation importance** -> knapsack problem.

End-to-end Cascades: Selecting Features

- Goal: Select feature set S that minimizes query time:

$$\min_S h_s \cdot \text{cost}(S) + (1 - h_s) \cdot \text{cost}(F)$$

- Approach:
 - Choose several potential values of $\text{cost}(S)$.
 - Find best feature set with each $\text{cost}(S)$.
 - **Train model & find cascade threshold for each set.**
 - Pick best overall.

End-to-end Cascades: Selecting Features

- Subgoal: Train model & find cascade threshold for S .

$$\min_S h_s \cdot \text{cost}(S) + (1 - h_s) \cdot \text{cost}(F)$$

- Solution:

- Compute empirically on held-out data.

End-to-end Cascades: Selecting Features

- Subgoal: Train model & find cascade threshold for S .

$$\min_S h_s \cdot \text{cost}(S) + (1 - h_s) \cdot \text{cost}(F)$$

- Solution:

- Compute empirically on held-out data.
- Train approximate model from S .

End-to-end Cascades: Selecting Features

- Subgoal: Train model & find cascade threshold for S .

$$\min_S h_s \cdot \text{cost}(S) + (1 - h_s) \cdot \text{cost}(F)$$

- Solution:
 - Compute empirically on held-out data.
 - Train approximate model from S .
 - Predict held-out set, determine cascade threshold empirically using accuracy target.

End-to-end Cascades: Selecting Features

- Goal: Select feature set S that minimizes query time:

$$\min_S h_s \cdot \text{cost}(S) + (1 - h_s) \cdot \text{cost}(F)$$

- Approach:
 - Choose several potential values of $\text{cost}(S)$.
 - Find best feature set with each $\text{cost}(S)$.
 - Train model & find cascade threshold for each set.
 - **Pick best overall.**

End-to-end Cascades: Results

- Speedups of up to 5x without statistically significant accuracy loss.
- Full evaluation at the end.

Outline

- System Overview
- Optimization 1: End-to-end Cascades
- **Optimization 2: Top-K Query Approximation**
- Evaluation

Top-K Approximation: Query Overview

- Top-K problem: Rank K highest-scoring items of a dataset.
- Top-K example: Find 10 artists a user would like most (recommender system).

Top-K Approximation: Asymmetry

- High-value items must be predicted, ranked precisely.
- Low-value items need only be identified as low value.

Artist	Score	Rank
Beatles	9.7	1
Bruce Springsteen	9.5	2
...
Justin Bieber	5.6	999
Nickelback	4.1	1000

High-value:
Rank precisely,
return.

Low-value:
Approximate,
discard.

Top-K Approximation: How it Works

- Use approximate model to identify and discard low-value items.
- Rank high-value items with powerful model.

Top-K Approximation: Prior Work

- Existing systems have similar ideas.
- However, we automatically generate approximate models for any ML application—prior systems don't.
- Similar challenges as in cascades.

Source: Cheng et al. (DLRS '16)

Top-K Approximation: Automatic Tuning

- Automatically selects features, tunes parameters to maximize performance given accuracy target.
- Works similarly to cascades.

Top-K Approximation: Results

- Speedups of up to 10x for top-K queries.
- Full evaluation at the end.

Outline

- System Overview
- Optimization 1: End-to-end Cascades
- Optimization 2: Top-K Query Approximation
- Evaluation

Willump Evaluation

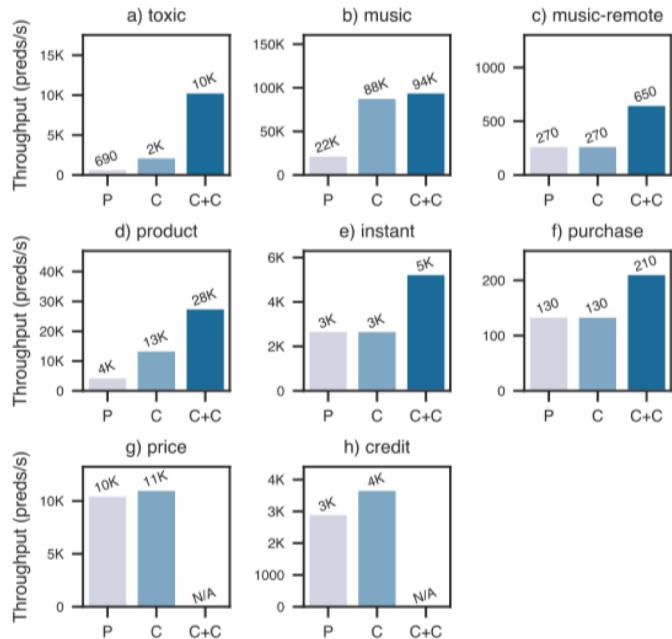


Figure 6. WILLUMP performance on offline batch queries. P means (unoptimized) Python, C means compiler optimizations only, C+C means compiler and cascades optimizations.

Willump Evaluation

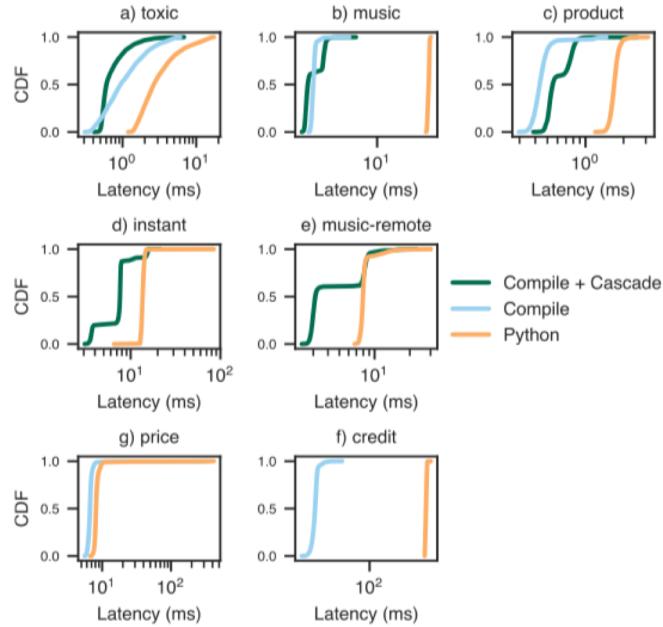


Figure 7. WILLUMP latency CDFs on online point queries, with one outstanding query at a time. Benchmarks in the second row contain no compilable operators; we only apply cascades. Benchmarks in the third row do not perform classification; we only apply compiler optimizations.

Willump Evaluation

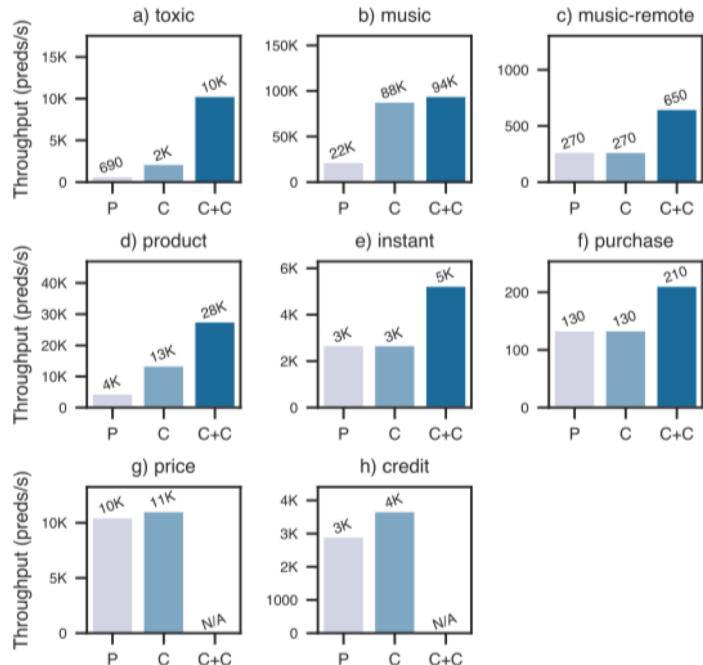


Figure 6. WILLUMP performance on offline batch queries. P means (unoptimized) Python, C means compiler optimizations only, C+C means compiler and cascades optimizations.

Willump Evaluation

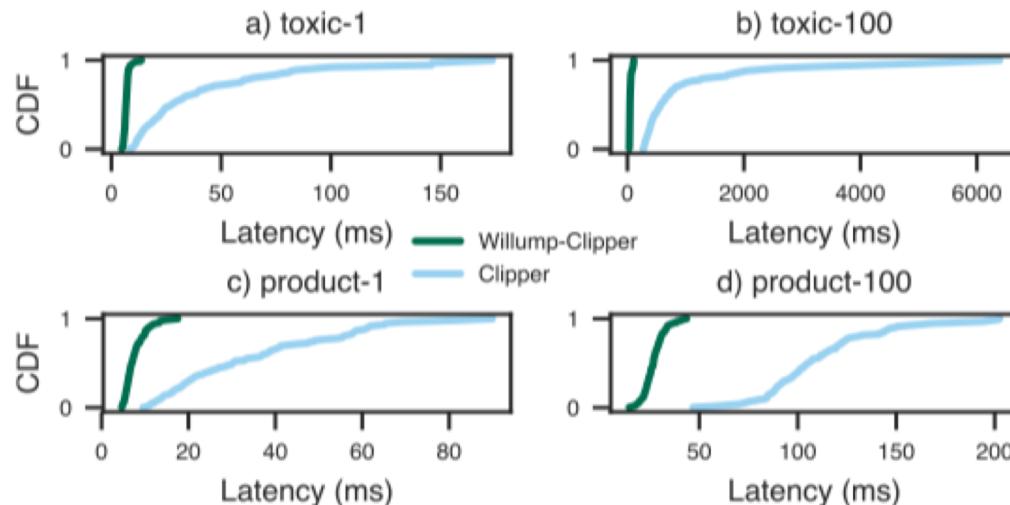


Figure 9. End-to-end latency CDFs on queries made using Clipper with and without WILLUMP optimization at batch sizes of 1 and 100.

Willump Evaluation

	Toxic	Music	Product	Instant	Purchase	Price	Credit
Cascades	1.1×	6.8×	2.2×	1.0×	1.7×	N/A	N/A
Top-K	1.1×	6.9×	1.8×	1.0×	1.6×	3.8×	3.9×

Table 2. Ratios of WILLUMP to Python training times when training end-to-end cascades and top-K query approximations.