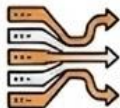# Bitcoin Fee Estimation:

## A Structural Model Approach

Kristian Praizner
Mentors: Dan Aronoff, Armin Sabouri

# Talk outline

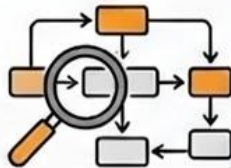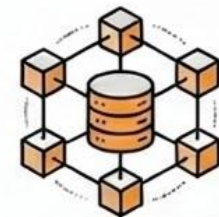- **Motivation**
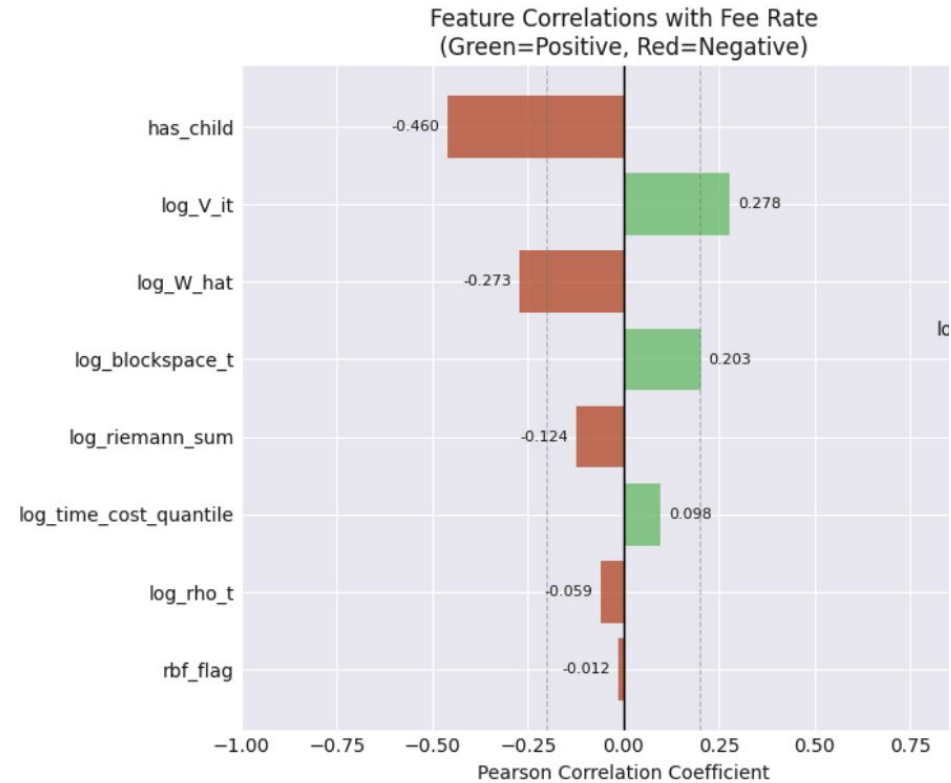- **Data Sources**
- **Methodology**
- **Conclusion**

Executive Summary:

- Main drivers of fee rate are V (transaction amount) + and CPFP (child pays for parent) -

- Several other variables are correlated with fee rate

- We choose to use a structural model in order to recover the true drivers by eliminating confounders



Feature Correlations with Fee Rate
(Green=Positive, Red=Negative)

| Feature | Pearson Correlation Coefficient |
|---|---|
| has_child | -0.460 |
| log_V_it | 0.278 |
| log_W_hat | -0.273 |
| log_blockspace_t | 0.203 |
| log_riemann_sum | -0.124 |
| log_time_cost_quantile | 0.098 |
| log_rho_t | -0.059 |
| rbf_flag | -0.012 |

# Executive Summary:

Main drivers of fee rate are V (transaction amount) + and CPFP (child pays for parent) -
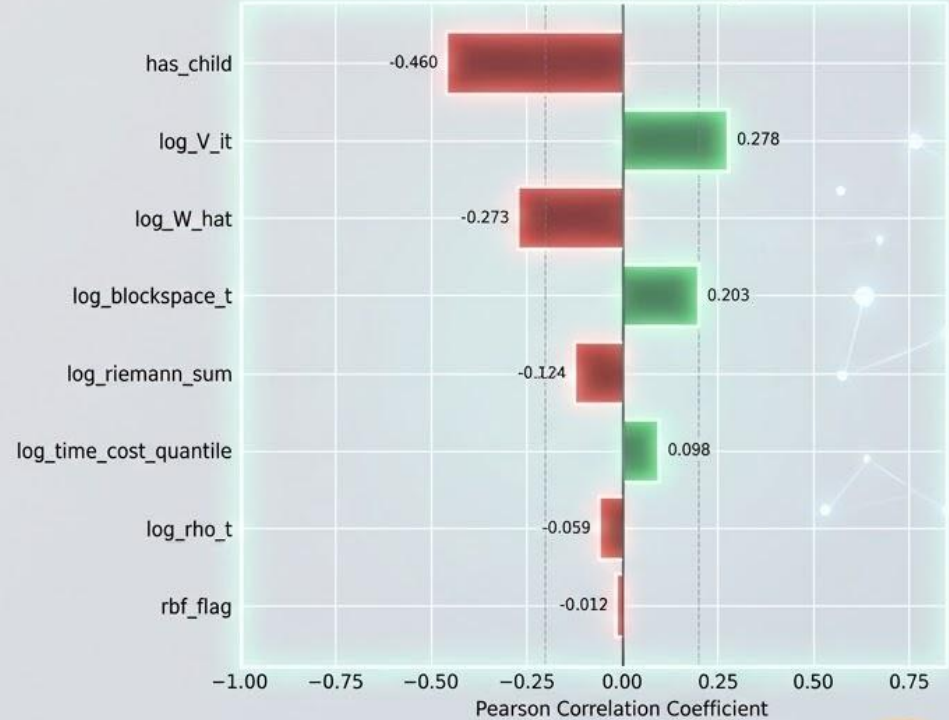
Several other variables are correlated with fee rate

We choose to use a structural model in order to recover the true drivers by eliminating confounders



Feature Correlations with Fee Rate
(Green=Positive, Red=Negative)

| Feature | Pearson Correlation Coefficient |
|---|---|
| has_child | -0.460 |
| log_V_it | 0.278 |
| log_W_hat | -0.273 |
| log_blockspace_t | 0.203 |
| log_riemann_sum | -0.124 |
| log_time_cost_quantile | 0.098 |
| log_rho_t | -0.059 |
| rbf_flag | -0.012 |

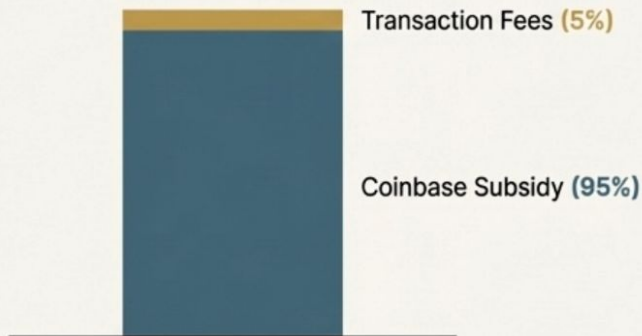# The Block Reward Has Two Components, But One is Programmed to Disappear.

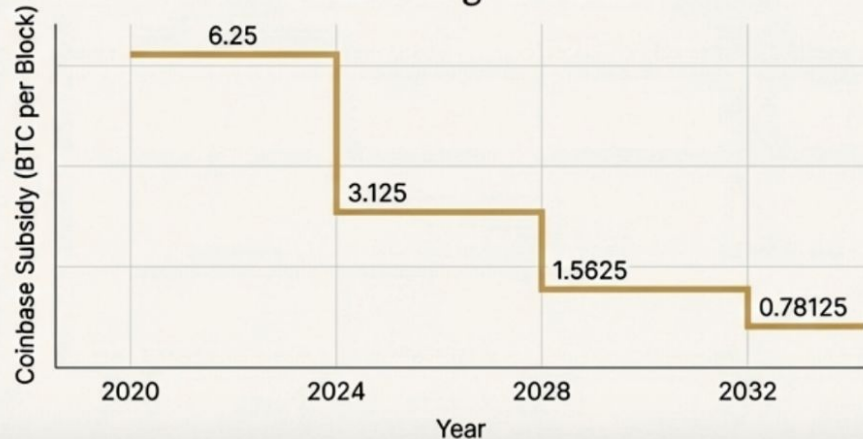The reward consists of two sources:
1. **Coinbase Subsidy**: A fixed amount of newly minted BTC, set by the protocol.
2. **Transaction Fees**: Voluntary fees paid by users to have their transactions included in a block.

Historically, the coinbase subsidy has comprised the vast majority of the reward. From 2022-2024, it accounted for approximately 95%.
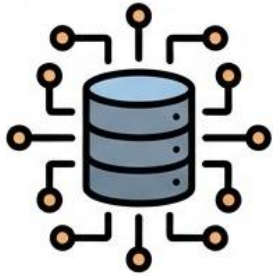
## Block Reward Composition (2022-2024)

Transaction Fees (5%)

Coinbase Subsidy (95%)

## The Halving Schedule

# To Maintain Security, Bitcoin's Price Must Grow Exponentially—An Unsustainable Path

**The Required Price Path to Offset the Halving**



If the BTC-denominated coinbase subsidy halves every four years, the USD price of Bitcoin must double in the same period just to keep the miners' revenue (and thus network security) constant. This implies a required annual growth rate of ~25%, a rate that exceeds historical trends and is dynamically impossible for any asset to sustain in the long term.

This raises the critical question for Bitcoin's future viability: **Will transaction fees rise to fill the ever-widening gap?**

**The model is Estimated on Granular, High-Frequency data from a Dedicated Bitcoin Node**

## Data Source & Structure

We operated a custom Bitcoin Node from August to December 2025 to collect high-fidelity mempool and blockchain data

## Data Structure

- The timeline was partitioned into epochs of 30 minutes
- For each transaction, we measured fee, mempool density, waittime, UTXO value, re-spend time, etc

## Implementation Details

The entire pipeline is available on Github

## Data Sourced

tx_id, tx_data, child_txid, conf_block_hash, found_at, mined_at, rbf_fee_total, min_respend_blocks, absolute_fee, fee_rate, version, seen_in_mempool, waittime, weight, size, total_output_amount, mempool_size, mempool_tx_count, output_weights

# The model is Estimated on Granular, High-Frequency data from a Dedicated Bitcoin Node

## Data Source & Structure

We operated a custom Bitcoin Node from August to December 2025 to collect high-fidelity mempool and blockchain data

## Data Structure 🕐

- The timeline was partitioned into epochs of 30 minutes
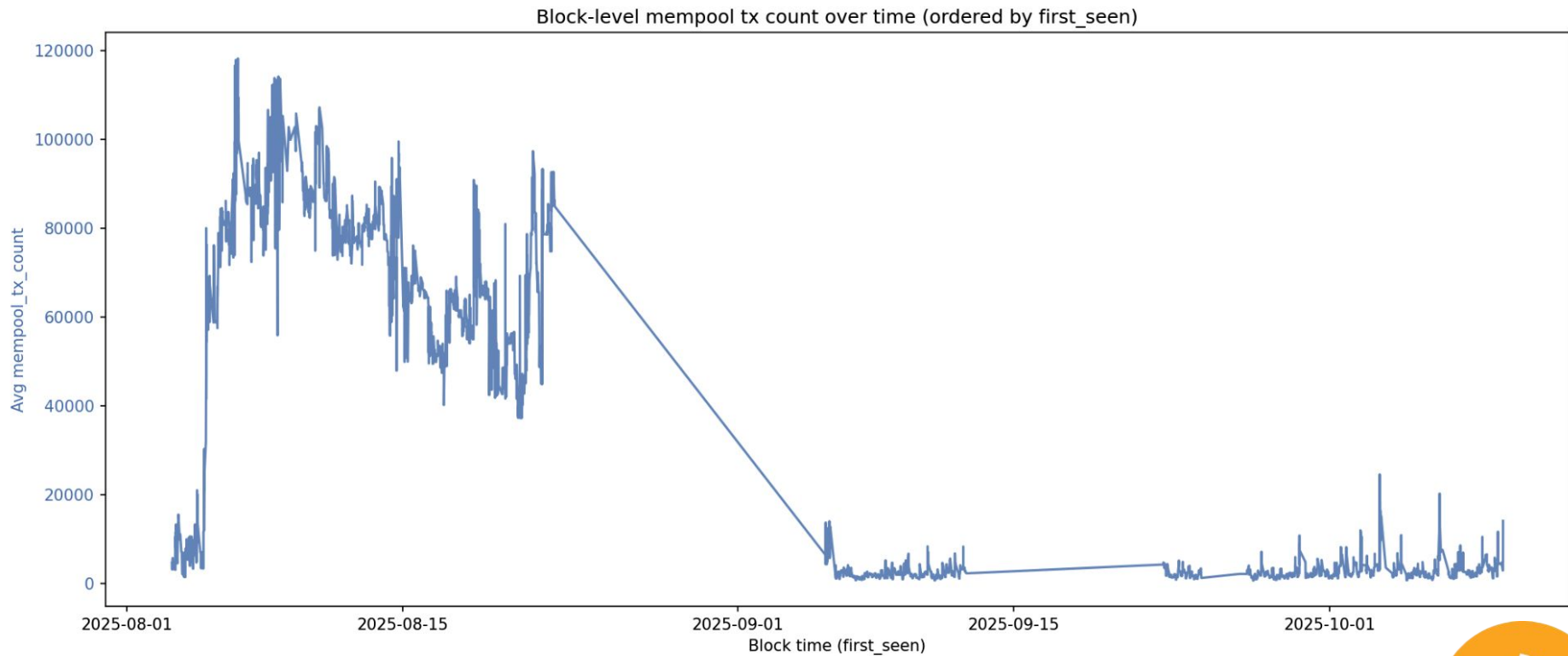- For each transaction, we measured fee, mempool density, waittime, UTXO value, re-spend time, etc

## Implementation Details

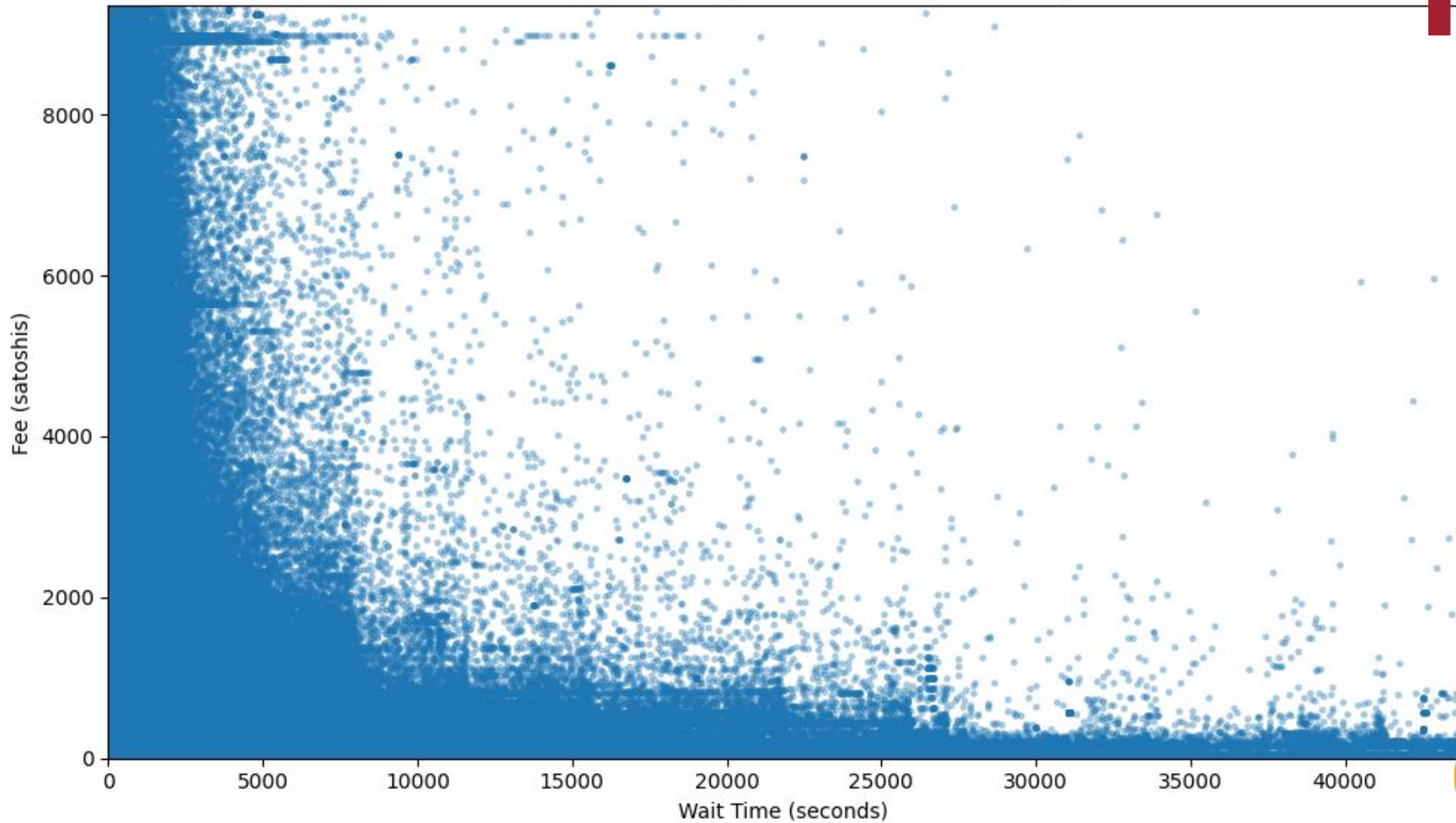The entire pipeline is available on Github

## Data Sourced

tx_id, tx_data, child_txtid, conf_block_hash, found_at, mined_at, rbf_fee_total, min_respend_blocks, absolute_fee, fee_rate, version, seen_in_mempool, waittime, weight, size, total_output_amount, mempool_size, mempool_tx_count, output_weights
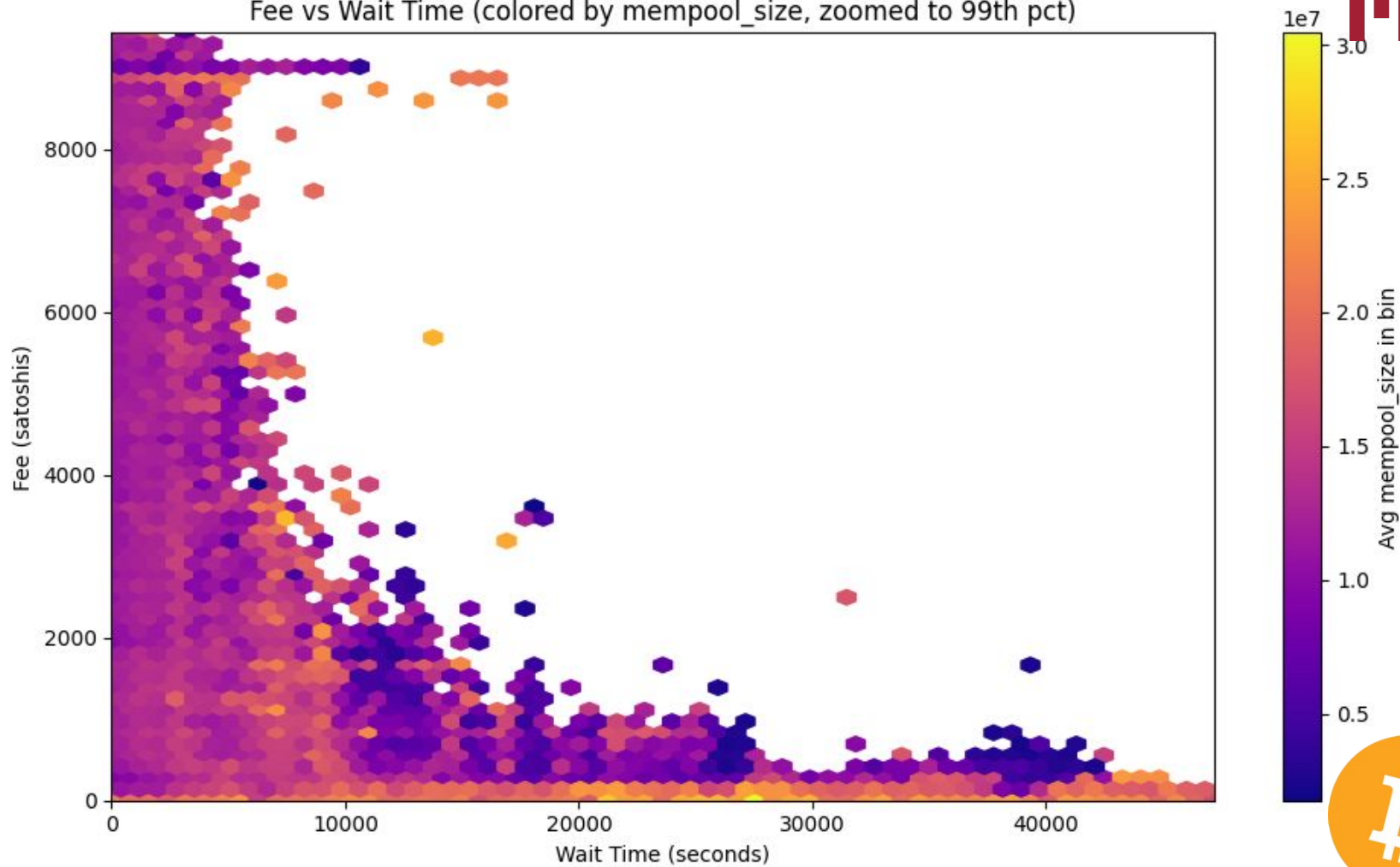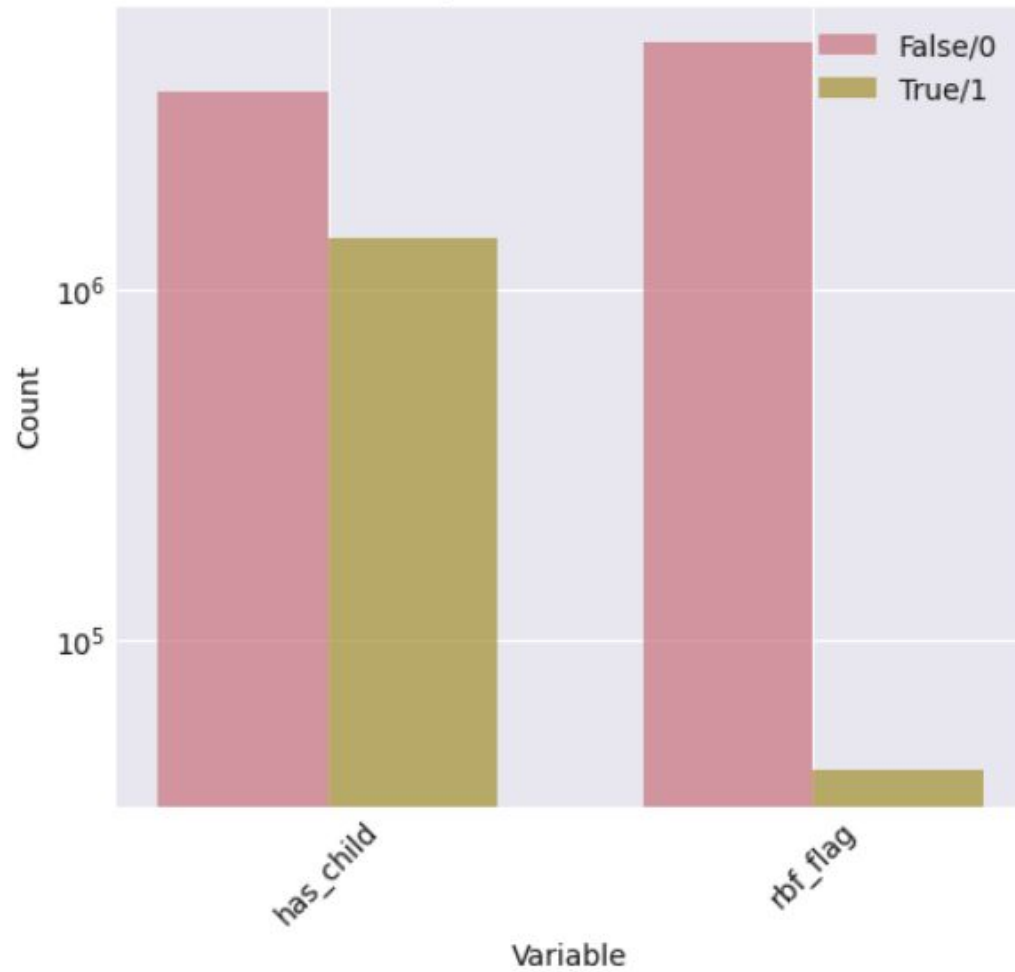
Block-level mempool tx count over time (ordered by first_seen)
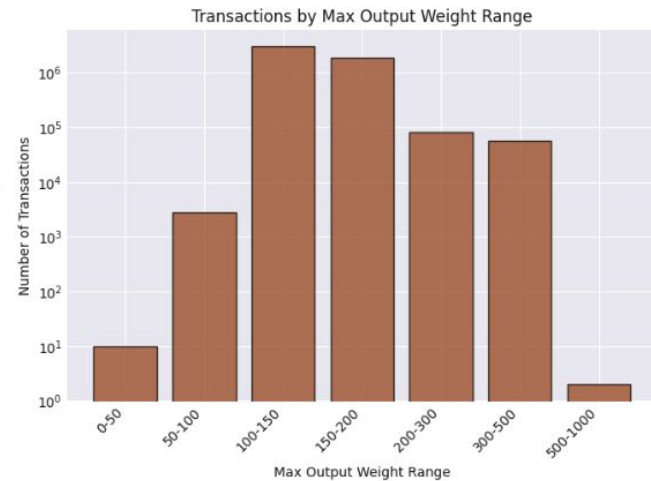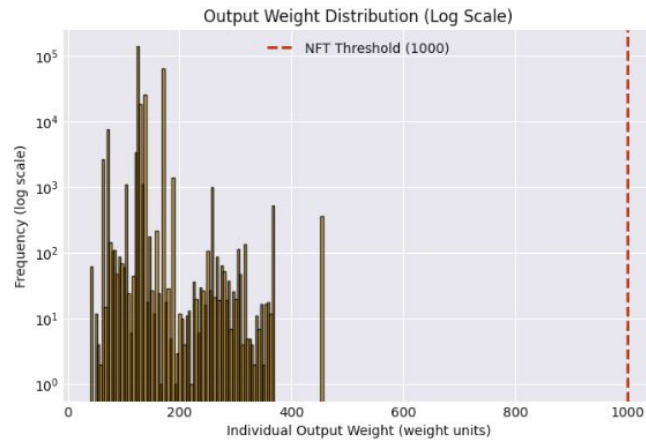
Fee vs Wait Time (zoomed to 99th percentile)

Fee vs Wait Time (colored by mempool_size, zoomed to 99th pct)

Binary Indicator Variables

## Block Fullness (by confirmed block)

## Blocks Above Utilization Thresholds

Distribution of Individual Output Weights
(sampled 267,398 outputs)

Distribution of Max Output Weight (per transaction)

Output Weight Distribution (Log Scale)

Transactions by Max Output Weight Range

# Methodology

Implementation of a two-stage structural model.

# We Model the Fee Market as a Queue Where Impatient Users Pay to Cut the Line.

- Our approach is built on the economic model of Huberman et al. (2021).

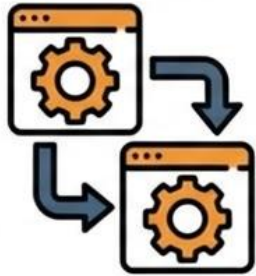- Miners, as profit-maximizers, prioritize transactions with the highest fees.

- Transactors differ in their "impatience" or time-cost. They select a fee to secure a desired spot in the queue of pending transactions (the mempool).

- The resulting fee is therefore a function of network congestion and the distribution of impatience among all users.

# The Challenge: How Can We Empirically Measure a User's 'Impatience'?

Transactor's
Impatience (c)

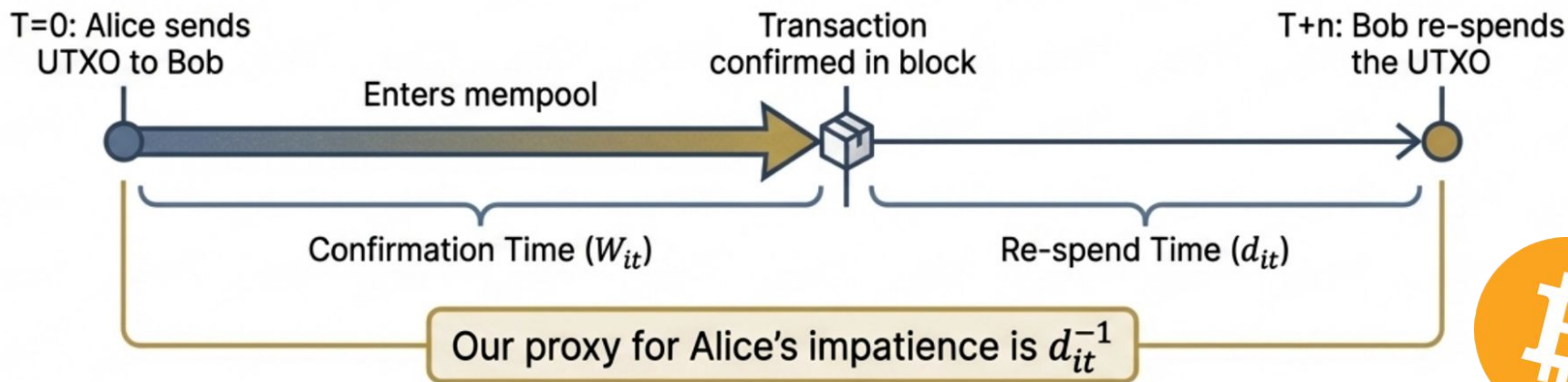Estimating the model requires a data-driven measure of impatience. This presents a core difficulty, as impatience (a user's "time-cost") is an unobservable internal state of mind. We need a clever, empirical proxy to extract this preference from on-chain data.

# The Breakthrough: Using 'Re-Spend Time' as a Proxy for Impatience

- We borrow an ingenious insight from Möser & Böhme (2015).
- We hypothesize that the time it takes for a transaction's output (UTXO) to be spent again is correlated with the original sender's urgency.
- A fast re-spend implies high impatience; a slow re-spend implies low impatience. The inverse of the re-spend time becomes our empirical proxy.

## Measuring Impatience via Re-Spend Time

T=0: Alice sends UTXO to Bob

Transaction confirmed in block

T+n: Bob re-spends the UTXO

Enters mempool

Confirmation Time ($W_{it}$)

Re-spend Time ($d_{it}$)

Our proxy for Alice's impatience is $d_{it}^{-1}$

Fee vs Wait Time (zoomed to 99th percentile)

Distribution of Blocks to Respend (0-20 blocks)

# A Random Forest Regressor for Wait Time Prediction

Our model uses a Random Forest Regressor to predict the log-transformed wait time (log_waittime) based on four key log-transformed and binary features.

## Model Overview:
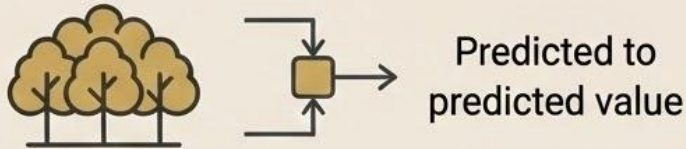
Predicted to predicted value

An ensemble learning method that constructs a multitude of decision trees at training time. It outputs the average prediction of the individual trees, reducing overfitting and improving accuracy.

## Features & Target (Log-Transformed):

**Target (y):** log_waittime (Log of Transaction Wait Time)

**Features (X):**
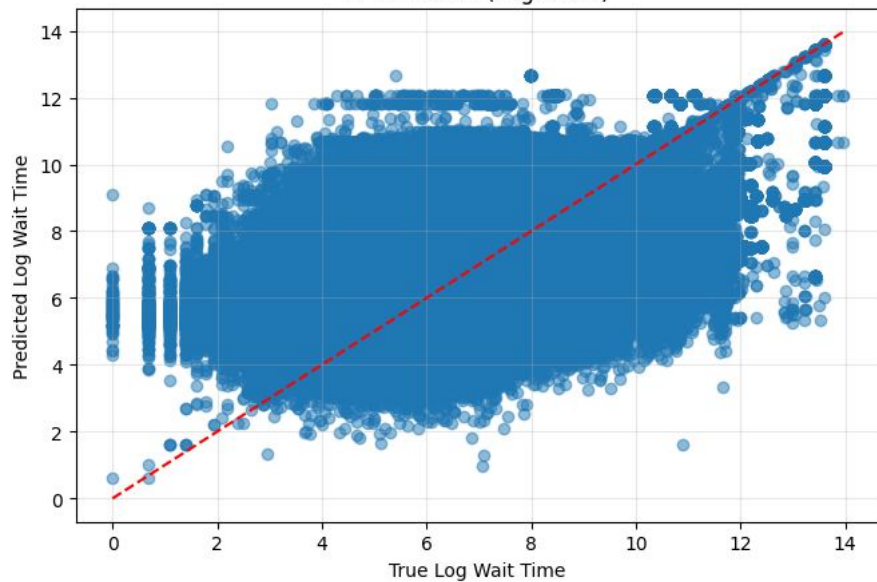- log_rho_t (Log of Network Congestion)
- log_time_cost (Log of Impatience Proxy)
- has_child (Binary: Has Child Transaction)
- rbf_flag (Binary: Replace-By-Fee Flag)

Random Forest Performance
R² = 0.4110 (Log Scale)

Feature Importance

**Actual vs Predicted Distribution (RF Model) - Extended Range**

**Confusion Matrix (Row-normalized %) - Extended Bins**

**Cumulative Distribution Function**

Rationale

- We use the first stage estimation to eliminate the influence of fees in waittime

- We do this by choosing variables that are not plausibly influenced by fees (congestion, time cost, CPFP)

- This proxy is then used as a variable in our second stage fee model

# Rationale

- We use the first stage estimation to eliminate the influence of fees in waittime

- We do this by choosing variables that are not plausibly influenced by fees (congestion, time cost, CPFP)

- This proxy is then used as a variable in our second stage fee model



Eliminate fee influence in waittime.

Choose variables not influenced by fees.

Proxy used as variable in final model.

# This Leads to Our Structural Model for Estimating Transaction Fees.

Our model empirically realizes the theory, expressing the transaction fee as a function of congestion, the aggregate impatience of other users, and other key transaction characteristics. The model is log-linear to reflect the "fat tail" distribution of fees.

**The Fee.** The transaction fee (in USD) that we aim to predict.

**Mempool Congestion.** The average number of transactions in the mempool during the epoch. How crowded is the waiting room?

$$b_{it} = \alpha_1 + \alpha_2 \hat{\rho}_t + \alpha_3 * \left[ \sum \ldots \right] + \alpha_4 V_{it} + \alpha_5 \text{Weight}_{it} + \ldots + \epsilon_{it}$$

**The Impatience Premium.** This core term aggregates the effect of all more-impatient users ahead in the queue, based on our re-spend time proxy.

**Control Variables.** We control for transaction value, weight (size), exchange activity, and other factors.

## Feature Correlations with Fee Rate
(Green=Positive, Red=Negative)

| Feature | Pearson Correlation Coefficient |
|---|---|
| has_child | -0.460 |
| log_V_it | 0.278 |
| log_W_hat | -0.273 |
| log_blockspace_t | 0.203 |
| log_riemann_sum | -0.124 |
| log_time_cost_quantile | 0.098 |
| log_rho_t | -0.059 |
| rbf_flag | -0.012 |

## Full Correlation Matrix
(All Model Features + Fee Rate)

| | log_rho_t | log_W_hat | log_riemann_sum | log_time_cost_quantile | log_V_it | has_child | rbf_flag | log_blockspace_t | log_fee_rate |
|---|---|---|---|---|---|---|---|---|---|
| log_rho_t | 1.00 | | | | | | | | |
| log_W_hat | 0.29 | 1.00 | | | | | | | |
| log_riemann_sum | 0.41 | 0.67 | 1.00 | | | | | | |
| log_time_cost_quantile | 0.14 | 0.01 | -0.19 | 1.00 | | | | | |
| log_V_it | -0.11 | -0.13 | -0.10 | -0.04 | 1.00 | | | | |
| has_child | 0.17 | 0.40 | 0.17 | -0.04 | -0.06 | 1.00 | | | |
| rbf_flag | 0.01 | 0.09 | 0.04 | -0.00 | 0.03 | -0.05 | 1.00 | | |
| log_blockspace_t | 0.12 | -0.17 | -0.08 | 0.02 | 0.06 | -0.13 | -0.06 | 1.00 | |
| log_fee_rate | -0.06 | -0.27 | -0.12 | 0.10 | 0.28 | -0.46 | -0.01 | 0.20 | 1.00 |

# Bitcoin Fee Estimation: A Structural Model Approach

## Model 1: Huber Robust Regression Analysis

### Model Performance Metrics

| Metric | Training Set | Test Set |
| --- | --- | --- |
| **R² (R-squared)** | 0.3067 | 0.3080 |
| **MAE (Mean Absolute Error)** | 0.44 | 0.44 |
| **Median AE (Median Absolute Error)** | 0.36 | 0.36 |
| **RMSE (Root Mean Squared Error)** | 0.59 | 0.59 |

Outliers Down-weighted
**1,291,475 (32.3%)**

### Top 10 Feature Coefficients (Standardized)

- Positive Correlation (Green)
- Negative Correlation (Red)

| Feature | Coefficient |
| --- | --- |
| has_child | -0.3054 |
| log_V_it | 0.1758 |
| log_time_cost_quantile | 0.0784 |
| log_blockspace_t | 0.0782 |
| log_W_hat | -0.0647 |
| log_riemann_sum | 0.0439 |
| log_weight | -0.0430 |
| rbf_flag | -0.0151 |
| log_rho_t | -0.0087 |
| V_it | -0.0008 |

Huber: Actual vs Predicted (Log Scale)

# Bitcoin Fee Estimation: A Structural Model Approach

## Model 2: Quantile Regression Analysis

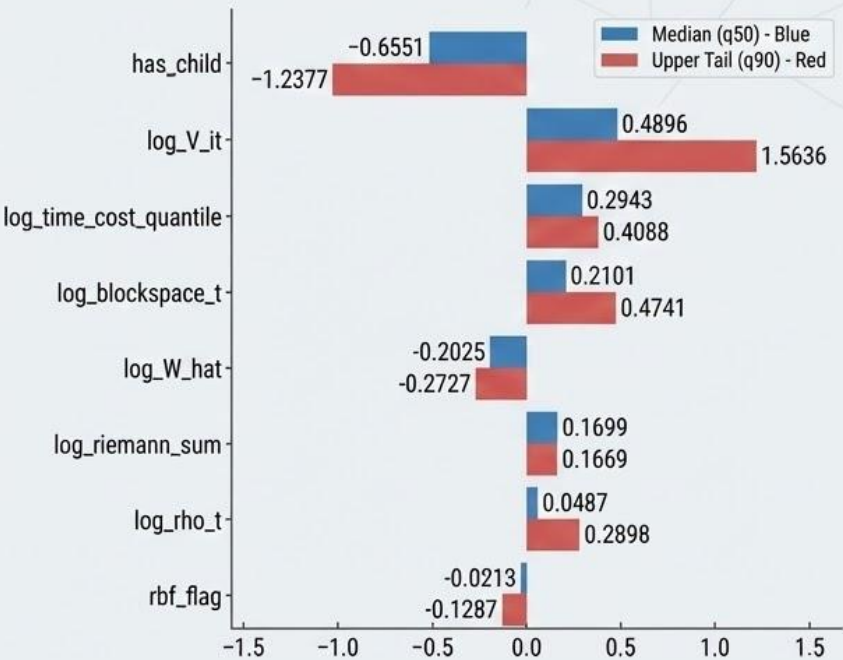## Model Performance Metrics

### Median Regression (50th percentile)

| Metric | Training Set | Test Set |
|---|---|---|
| Pseudo-$R^2$ | 0.1029 | 0.1038 |
| MAE | 1.93 | 1.96 |
| Median AE | 0.91 | 0.90 |

### 90th Percentile Regression (upper tail)

| Metric | Test Set (90th percentile) |
|---|---|
| Pseudo-$R^2$ | 0.0793 |
| MAE | 4.14 |

Using 10,000 samples for quantile regression (computational efficiency)

## Coefficient Comparison: Median vs 90th Percentile



Legend:
- Median (q50) - Blue
- Upper Tail (q90) - Red

| Variable | Median (q50) | Upper Tail (q90) |
|---|---|---|
| has_child | −0.6551 | −1.2377 |
| log_V_it | 0.4896 | 1.5636 |
| log_time_cost_quantile | 0.2943 | 0.4088 |
| log_blockspace_t | 0.2101 | 0.4741 |
| log_W_hat | −0.2025 | −0.2727 |
| log_riemann_sum | 0.1699 | 0.1669 |
| log_rho_t | 0.0487 | 0.2898 |
| rbf_flag | −0.0213 | −0.1287 |

Model 2: Quantile Regression Diagnostics (q50 & q90)

# Bitcoin Fee Estimation: A Structural Model Approach
## Model 3: Spline Regression (Segmented)

## Model Configuration & Features

📊 Using: 50,000 samples for Spline Regression

📈 Spline features (non-linear): ['log_rho_t', 'log_V_it', 'log_W_hat', 'log_blockspace_t']

📉 Linear features: ['log_weight', 'log_riemann_sum', 'log_time_cost_quantile', 'has_child', 'rbf_flag', 'V_it']

⚙️ Spline configuration:
- Knots: 30 (creates 29 segments)
- Degree: 3 (cubic splines)
- Features per spline variable: 32

➡️ Fitting spline regression model (log1p target)...

## Model Performance Metrics

| Metric | | Training Set | Test Set |
|---|---|---|---|
| $R^2$ | 🔴 | -0.0298 | -0.0561 |
| MAE | 🟢 | 2.99 | 2.99 |
| Median AE | 🟢 | 1.89 | 1.89 |
| RMSE | 🟢 | 5.52 | 5.47 |

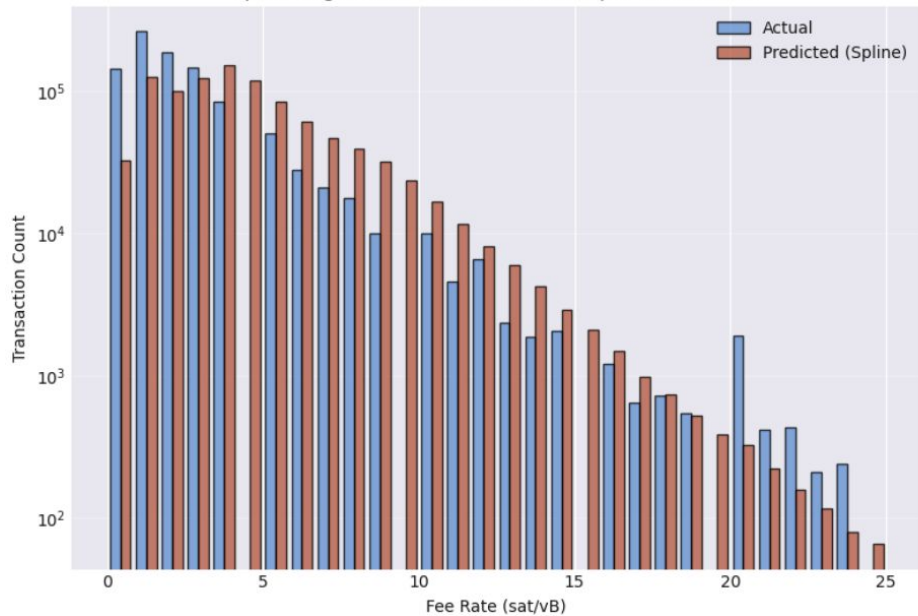## Feature Transformation Summary

Total features after spline transformation: 130
- Spline-derived features: 124
- Linear features: 6

## Spline Feature Importance (sum of |coefficients|)



- splines_block: 16.8912 Spline
- has_child: 0.2379 Linear
- V_it: 0.1308 Linear
- log_weight: 0.0942 Linear
- rbf_flag: 0.0374 Linear
- log_riemann_sum: 0.0273 Linear
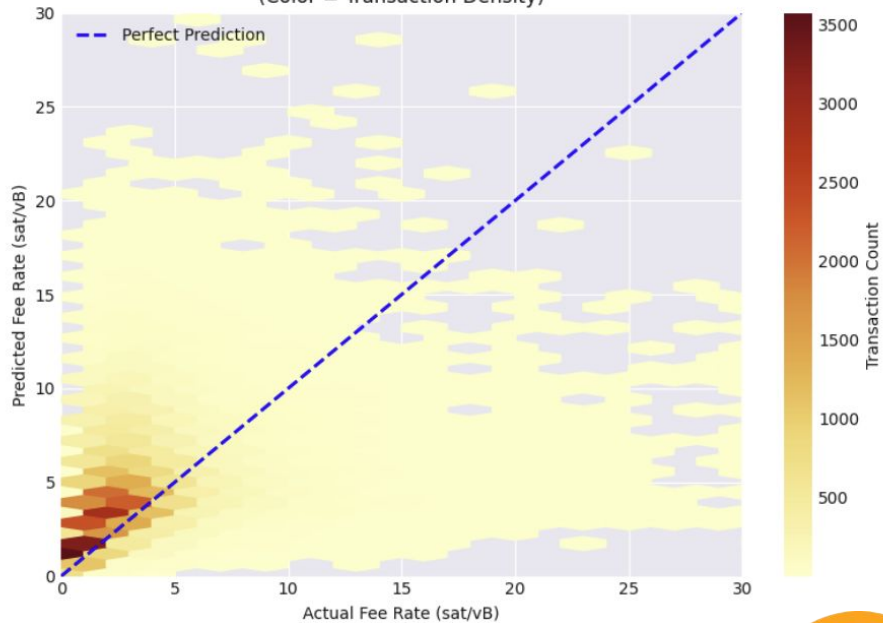- log_time_cost_quantile: 0.0237 Linear

Spline Regression: Transaction Count per Fee Rate Bin

Spline Regression: Actual vs Predicted
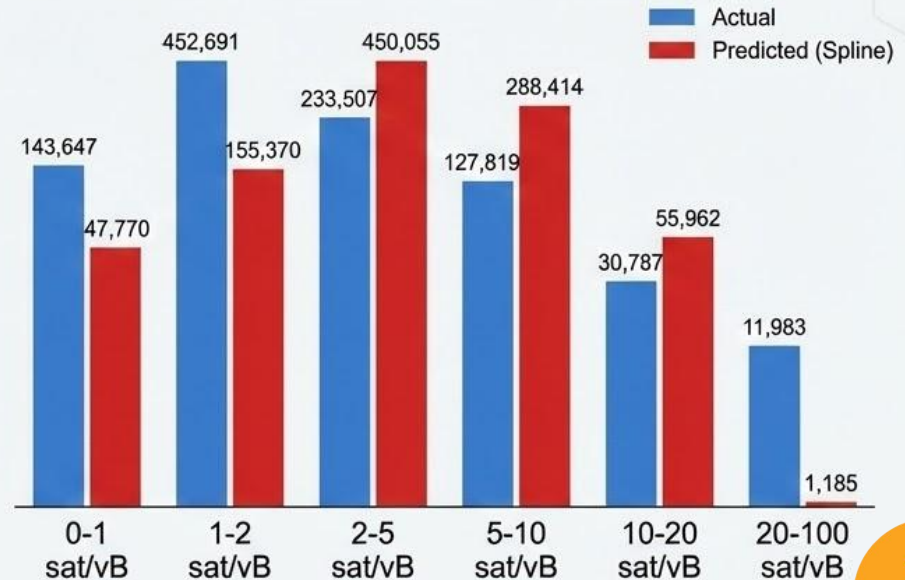(Color = Transaction Density)

# Bitcoin Fee Estimation: A Structural Model Approach
## Spline Regression: Fee Rate Distribution Summary

## Fee Rate Distribution Statistics

| Metric | Actual | Predicted (Spline) |
|---|---|---|
| **Total Transactions:** | 1,000,789 | - |
| **Mean** (sat/vB): | 3.09 | 4.56 |
| **Median** (sat/vB): | 2.00 | 3.96 |
| **Std Dev** (sat/vB): | 5.32 | 3.02 |

## Transaction Count by Fee Rate Bracket



Legend: Actual (blue), Predicted (Spline) (red)

| Bracket | Actual | Predicted (Spline) |
|---|---|---|
| 0-1 sat/vB | 143,647 | 47,770 |
| 1-2 sat/vB | 452,691 | 155,370 |
| 2-5 sat/vB | 233,507 | 450,055 |
| 5-10 sat/vB | 127,819 | 288,414 |
| 10-20 sat/vB | 30,787 | 55,962 |
| 20-100 sat/vB | 11,983 | 1,185 |

# Bitcoin Fee Prediction: Scenario Analysis

Based on Input Parameters (Model 3: Spline Regression)

## Scenario 1: Lower Mempool Density

**INPUTS**

| | | |
|---|---|---|
| rho_t | 8,000 | |
| blockspace_t | 0.3500 | |
| V_it | 200,000 | |
| has_child | No (0) | ✗ |
| rbf_flag | Yes (1) | ✓ |

**PREDICTED FEE (sat/vB)**

### 2.0992

pred_fee_spline_sat_vB

## Scenario 2: Higher Mempool Density

**INPUTS**

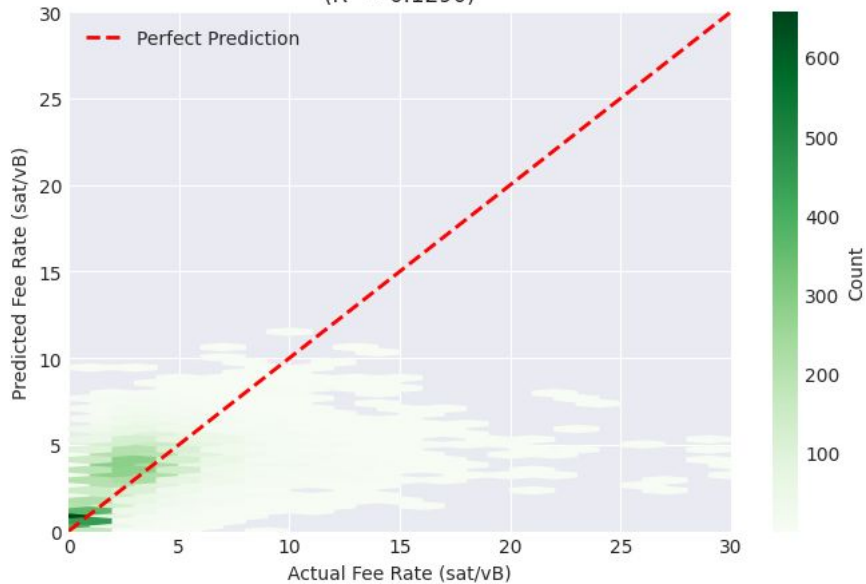| | | |
|---|---|---|
| rho_t | 60,000 | |
| blockspace_t | 0.8500 | |
| V_it | 5,000,000 | |
| has_child | Yes (1) | ✓ |
| rbf_flag | No (0) | ✗ |

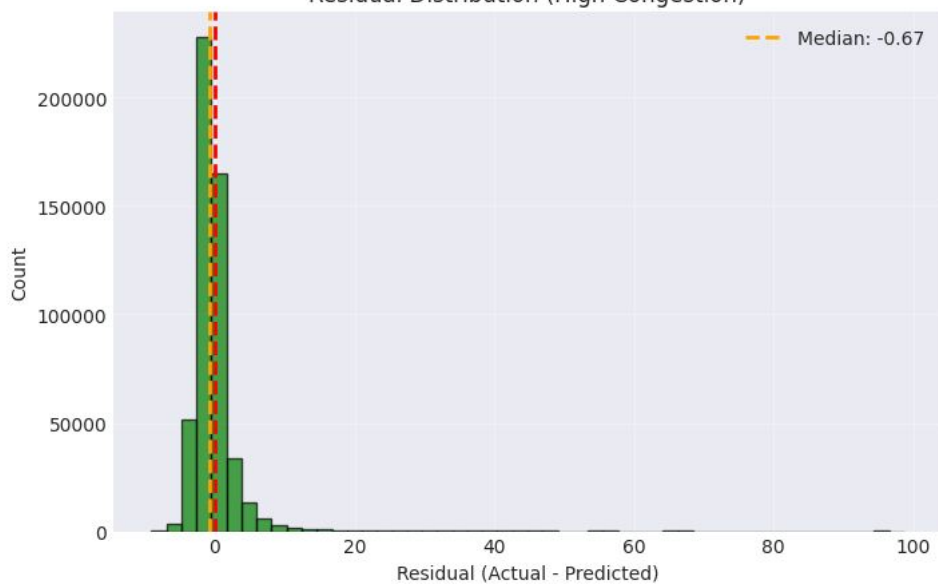**PREDICTED FEE (sat/vB)**

### 1.2015

pred_fee_spline_sat_vB

High Congestion: Actual vs Predicted
(R² = 0.1290)

Residual Distribution (High Congestion)

Median: -0.67

# SUMMARY: F-STATISTICS

## LINEAR MODEL F-STATISTICS

log_V_it: 19505.97 (***)

has_child: 16780.39 (***)

log_rho_t: 5625.70 (***)

log_blockspace_t: 1173.33 (***)

log_time_cost_quantile: 657.31 (***)

log_W_hat: 580.18 (***)

rbf_flag: 104.35 (***)

log_riemann_sum: 89.73 (***)

## NONLINEARITY?

↗ YES (R2 Gain: 0.13%)

N/A (binary)

↗ YES (R2 Gain: 0.66%)

↗ YES (R2 Gain: 0.17%)

↗ YES (R2 Gain: 0.02%)

↗ YES (R2 Gain: 0.15%)

N/A (binary)

↗ YES (R2 Gain: 0.05%)

## SPLINE MODEL F-STATISTICS

log_V_it: 95.70 (***)

has_child: N/A

log_rho_t: 475.86 (***)

log_blockspace_t: 123.70 (***)

log_time_cost_quantile: 12.90 (***)

log_W_hat: 106.83 (***)
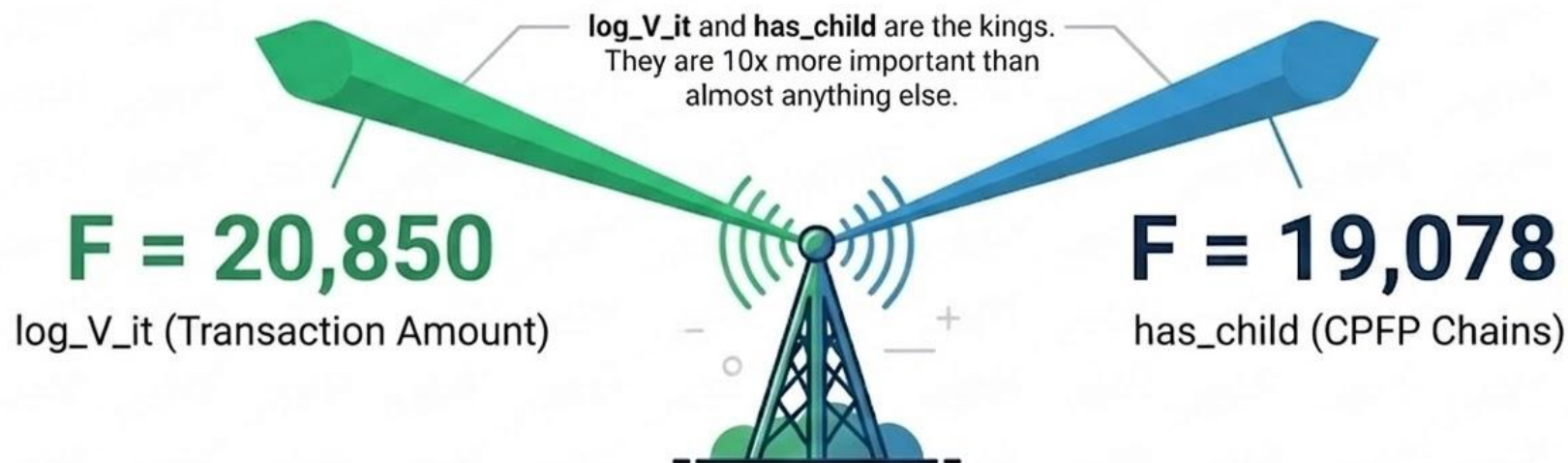
rbf_flag: N/A

log_riemann_sum: 33.08 (***)

*** $p < 0.001$

Summary and Conclusions

# THE GOOD: FOUND THE SIGNAL

F-statistics are **massive for the top features**. This confirms that these **variables are definitely drivers of fee rates.**

**log_V_it** and **has_child** are the kings. They are 10x more important than almost anything else.

F = 20,850

log_V_it (Transaction Amount)

F = 19,078

has_child (CPFP Chains)

**THE TAKEAWAY:**
Transaction Amount and CPFP (Child-Pays-For-Parent) chains are the primary drivers of fees.

# KEY TAKEAWAYS:



Modeling BTC fees is hard (power law)

Machine Learning Models would be better

Mempool congestion hardly affects fee rate at all

# ACKNOWLEDGMENTS