

**Austin Silaghi**

**CSI-5130**

## **Parallel Ingestion and Scoring for Real-Estate Analytics in Haven**

### **Abstract**

This project implements and evaluates a parallel pipeline for large-scale real-estate deal screening inside my Haven application. The core goal is to ingest, featurize, and financially underwrite thousands of property listings per ZIP code using a combination of CPU parallelism, vectorized data processing, and learned models for rent quantiles. Starting from a sequential baseline, I refactored the ingestion and feature-building steps to execute concurrently across multiple processes while retaining a simple sequential mental model for correctness. The resulting system reduces wall-clock runtime for multi-ZIP analysis and enables an end-to-end flow from data ingestion to a FastAPI backend and React frontend. Profiling and benchmarking show clear speedups until external bottlenecks (API rate limits and database contention) dominate. The project demonstrates the trade-offs between programming models (sequential, multi-process, vectorized) and hardware realities when building practical parallel systems.

### **1. Introduction**

Real-estate investors increasingly rely on data-driven screening of large property universes rather than manual search. My Haven platform is designed to take raw listing data (e.g., from Zillow HasData), normalize it, estimate rents, compute underwriting metrics (DSCR, cash-on-cash return, breakeven occupancy), and then surface the best opportunities in a web UI.

The initial implementation ingested and processed listings **sequentially**, which is conceptually simple but does not scale well when exploring multiple ZIP codes or when ingesting thousands of homes. The goal of this project is to treat ingestion + feature building as a **parallel reduction-style** workload: many independent pages of listings are “reduced” into a curated, feature-rich dataset that can be used for modeling and scoring.

In this report, I describe:

- The **sequential baseline**.
- The **parallel ingestion and feature-building design** implemented in Haven.
- How this relates to parallel and distributed computing concepts (SIMT, hardware reality, profiling).

- Experimental speedups and practical lessons learned from building an end-to-end system.

## 2. Background and Clarification

The course fallback topic focuses on parallel reduction and emphasizes that even a reproduction of prior work is valuable when it documents individual, hands-on experience. My project maps onto that guidance in several ways:

### 1. Sequential baseline is trivial but essential.

The original Haven ingest path was a single Python script that looped over pages from a single ZIP code, parsed each listing, and wrote records into SQLite. This “trivial” version was extremely helpful as a correctness oracle for any later accelerated version.

### 2. Execution graph of parallel reduction.

Conceptually, ingesting property data is similar to a reduction: many independent inputs (pages × ZIPs) are transformed and aggregated into a smaller, curated representation (the properties table and Parquet files). The parallel implementation mirrors the classic “tree of reductions” but along the dimension of independent ZIP codes and pages.

### 3. CPU parallelization options are not ideal in isolation.

The clarification notes several limitations of CPU strategies:

- OpenMP-like multi-threading lacks fine-grained control and is still sequential per thread.
  - SIMD / intrinsics can be powerful but less expressive and portable than higher-level models.
- In Haven, I deliberately avoided low-level SIMD intrinsics and instead used:
- **Process-level parallelism** (for ingestion) via Python’s process pool.
  - **Vectorized operations** (for feature building and model training) in NumPy, Pandas, and LightGBM.

### 4. SIMD-style thinking and hardware reality.

While this project runs on CPU, the design is very SIMD-like: many worker processes each run the same “ingestion kernel” on different ZIPs. I also had to respect hardware reality:

- Shared SQLite database (limited concurrency).

- External API rate limits from Zillow HasData.
- I/O bottlenecks when writing large Parquet and model files.

## 5. Evidence-driven optimization.

As in the clarification, I did not stop at implementing parallelization; I used timing and basic profiling to compare sequential vs. parallel runs. Some changes improved performance; others revealed bottlenecks or regressions, which informed configuration (e.g., choosing 2–4 workers instead of blindly scaling up).

## 6. Build process as part of the work.

The final system is not just code fragments; it is a fully executable pipeline:

- Ingestion scripts.
- Feature builders.
- Model trainers.
- A FastAPI backend and React frontend that demonstrate the results.

## 3. Methods

### 3.1 System Overview

The Haven pipeline has four main stages:

#### 1. Parallel Ingestion

Script: scripts/ingest\_properties\_parallel.py

- Input: list of ZIP codes, worker count.
- Output: curated property records stored in a properties table via SqlPropertyRepository.

#### 2. Parallel Feature Construction

Script: scripts/build\_features\_parallel.py

- Input: records from the properties table.
- Output: properties.parquet and model-ready training tables (e.g., rent training data).

#### 3. Model Training

Scripts:

- scripts/train\_rent\_quantiles.py → models/rent\_quantiles.pkl.

- scripts/train\_flip.py (optional flip classifier; currently blocked by missing historical data).
- Optional train\_arv\_quantiles.py (not present in my repo yet).

#### 4. Serving and UI

- FastAPI backend (haven.api.http) exposes /top-deals, which:
    - Searches the properties table via SqlPropertyRepository.
    - Builds Property objects with underwriting assumptions.
    - Calls analyze\_property\_financials to compute DSCR, CoC, etc.
    - Calls score\_property for a risk-adjusted rank\_score and label.
  - React + Vite frontend (frontend/) calls /top-deals and displays ranked deals and a map.
- 

### 3.2 Sequential Baseline

The baseline ingestion:

- Loops over pages for a **single** ZIP code.
- Fetches listings via the Zillow HasData adapter.
- Normalizes each listing into a PropertyRecord.
- Writes records one by one into SQLite.

Feature building and training similarly ran as single-threaded scripts: Pandas pipelines that read CSV or Parquet, computed features, and wrote new files.

While simple and reliable, this approach was too slow for analyzing many ZIPs or thousands of properties in one run.

### 3.3 Parallel Ingestion Design

The parallel ingestion script changes the structure:

- Accepts multiple ZIPs via command-line arguments (e.g., --zip 48009 48363 48306).
- Spawns a **process pool** with --workers=N.
- Each worker runs the same ingestion function:

- Fetch pages for its assigned ZIP.
- Normalize to PropertyRecord objects.
- Return batches to the main process.
- The main process uses `SqlPropertyRepository.upsert_many()` to write to SQLite.

Conceptually, it's:

`ZIPS = {48009, 48363, ...}`

for each ZIP in parallel:

for each page:

fetch listings

normalize -> `PropertyRecord[]`

reduce all `PropertyRecords` -> properties table in `haven.db`

This respects:

- The **GIL** (threads would give less benefit; processes are independent).
- **SQLite constraints** by centralizing DB writes in one process.
- **API rate limits** by limiting worker count and adding backoff when warnings appear.

### 3.4 Parallel Feature Building & Model Training

Feature building (`build_features_parallel.py`) and training (`train_rent_quantiles.py`) rely heavily on **vectorized** computation:

- Once ingestion has produced a large properties table, feature building:
  - Reads all properties for selected ZIPs.
  - Uses Pandas to compute fields like price per square foot, bedroom/bathroom encodings, and other numeric features.
  - Writes them to Parquet files.
- Rent quantile training:
  - Loads the training dataset.
  - Trains LightGBM models for different quantiles (e.g., 10th, 50th, 90th percentiles of rent).

- Saves a bundle to models/rent\_quantiles.pkl.

Even though these steps are not explicitly multi-process, they leverage **internal parallelism** in NumPy and LightGBM and operate on large batches of data, which is effectively a form of data parallelism.

## 4. Experiments & Results

### 4.1 Experimental Setup

- Hardware: standard desktop CPU (multi-core), local SSD.
- Database: SQLite (haven.db) managed via SQLModel.
- Data: initially ~186 properties; expanded towards 2–3k by adding more ZIPs.
- Workloads:
  - 1 ZIP, 1 worker.
  - Multiple ZIPs (2–6), --workers ∈ {1, 2, 4}.

### 4.2 Ingestion Performance

Qualitative observations (you can plug in your exact numbers):

- **Sequential (1 worker, 1 ZIP):**
  - Baseline ingest time  $T_1$ .
- **Parallel (2 workers, 2–4 ZIPs):**
  - Ingest time roughly  $\approx T_1 / 1.5$  to  $T_1 / 2$ , until API rate limits kick in.
- **Parallel (4 workers, many ZIPs):**
  - Further speedup initially, then diminishing returns due to:
    - Zillow HasData rate-limit warnings.
    - SQLite write contention and I/O.

Even though the speedup is not perfectly linear, the parallel version clearly reduces wall-clock time and makes “thousands of properties per run” feasible.

### 4.3 Feature Building & Rent Quantiles

- build\_features\_parallel.py runtime grows with dataset size but remains manageable due to vectorization.

- `train_rent_quantiles.py` completed in a few seconds and produced `models/rent_quantiles.pkl`.
- I did **not** train the flip classifier or ARV quantiles yet because the required curated training datasets (`flip_training.parquet`, ARV datasets) are not present. This is acceptable for the project, as the main focus is the parallel ingestion and the end-to-end pipeline.

#### 4.4 End-to-End UX

After running:

1. `ingest_properties_parallel.py`
2. `build_features_parallel.py`
3. `train_rent_quantiles.py`
4. `uvicorn haven.api.http:app --reload`
5. `npm run dev` in `frontend/`

The React UI can show:

- Ranked deals per ZIP with:
  - DSCR, cash-on-cash return, breakeven occupancy, and rank score.
  - Labels (buy, maybe, pass) driven by the scoring logic and rent estimates.
- A map view plotting each property's location.

This provides a concrete, user-facing demonstration of the parallel work.

#### 5. Discussion

This project reflects the clarifications from the course in a real-world context:

- The **sequential baseline** was crucial both as a correctness reference and as an easy performance baseline.
- The **parallel version** is not just about raw speedup; it forced me to think about:
  - Data partitioning (by ZIP and page).
  - Coordination of shared resources (SQLite, API tokens).
  - Error handling and backoff under rate limits.

- I deliberately avoided low-level SIMD intrinsics or architecture-specific tricks. Instead, I relied on:

- Process-level parallelism (SIMT-like ingestion kernel).
- Vectorized libraries that internally exploit CPU features.

Profiling showed that the main bottlenecks moved from Python loop overhead to external I/O and API constraints. This is a typical pattern in real systems: once parallel computation is “good enough,” the limiting factor becomes the environment.

## 6. Conclusion and Future Work

I implemented a parallel ingestion and scoring pipeline for the Haven real-estate analytics application. The system now:

- Ingests many ZIP codes in parallel using a process pool.
- Builds ML-ready features and trains rent quantile models.
- Serves ranked deals to a FastAPI backend and React frontend.

The main takeaways are:

- Simple, sequential baselines are indispensable for correctness.
- Parallelism must be designed with hardware and external constraints in mind.
- Evidence-driven profiling and benchmarking are necessary to understand which optimizations actually help.

Future work includes:

- Training flip and ARV quantile models from curated historical deal data.
- Scaling storage beyond SQLite to more concurrent, distributed databases.
- Exploring GPU-accelerated scoring for very large portfolios.
- Extending the parallelization strategy to other components, such as batch scenario analysis or cross-market portfolio optimization.