# CS129 Project - Plan, Results

## Part 1

### ▼ Phase 1: Setup & Data Understanding (Days 1-2)

### Environment Setup

- ○ Set up Python environment (Python 3.10+)
- ○ Install packages: `numpy`, `pandas`, `scikit-learn`, `matplotlib`, `scipy`
- ○ Create project folder structure:

```
project/
├── data/
│   ├── uci_hapt/
│   ├── marques/
│   └── sisfall/
├── notebooks/
├── src/
└── results/
```

- ○ **Download & Explore UCI HAPT**
- ○ Download UCI HAPT dataset: https://archive.ics.uci.edu/dataset/341

- ◯ Read the README — understand what files are included
- ◯ Locate the **raw inertial signals** (not the pre-computed 561 features)
  - ☐ These are the accelerometer (x,y,z) and gyroscope (x,y,z) files
  - ☐ Sampled at 50Hz
- ◯ Locate the **labels file** — this tells you which time ranges correspond to which activity
- ◯ Locate the **subject IDs** — you need to know which data belongs to which of the 30 subjects
- ◯ Understand the 12 activity labels:
  - ☐ 6 static/dynamic: walking, walking upstairs, walking downstairs, sitting, standing, laying
  - ☐ 6 transitions: stand-to-sit, **sit-to-stand**, sit-to-lie, lie-to-sit, stand-to-lie, lie-to-stand
  - ☐ **You only care about sit-to-stand (label 5 in UCI HAPT)**
- ◯ Count how many sit-to-stand instances exist in the dataset

## Download External Datasets (for later, but grab now)

- ◯ Download Marques dataset (search: "Accelerometer data sit-to-stand elderly Marques 2020")
  - • 40 elderly participants, accelerometer only, waist smartphone
- ◯ Download SisFall dataset (search: "SisFall dataset")
  - • 15 elderly participants, accelerometer+gyroscope, waist device
- ◯ For each: skim the documentation, note the sampling rate, sensor types, and label format
  - • If sampling rate differs from 50Hz, you will need to resample using `scipy.signal.resample_poly` with anti-aliasing to avoid corrupted signals (especially SisFall at 200Hz → 50Hz)

▼ **Results**

Here's the summary of the UCI HAPT dataset analysis:

## Activity Labels

- ○ **SIT_TO_STAND = Activity ID 8** (confirmed)

## Accelerometer File Format

- 3 columns of float values (x, y, z acceleration), space-separated
- Files named `acc_expXX_userYY.txt` (one per experiment per subject)

## Labels File Format

Each row: `experiment_id user_id activity_id start_sample end_sample`

## Sit-to-Stand (Activity 8) Counts

- **Total instances: 62** across all 30 subjects
- **Per subject: 2 instances each**, except users 8, 22 have **3 instances** (due to extra experiments)
- Average duration: **~130 samples = ~2.6 seconds** at 50Hz (range: 74–183 samples)

## Class Imbalance

| Activity | Count |
| --- | --- |
| WALKING_DOWNSTAIRS (3) | 186 |
| WALKING_UPSTAIRS (2) | 183 |
| WALKING (1) | 127 |
| LAYING (6) | 120 |
| STANDING (5) | 120 |
| SITTING (4) | 120 |
| **SIT_TO_STAND (8)** | **62** |
| SIT_TO_LIE (9) | 60 |
| STAND_TO_SIT (7) | 60 |
| LIE_TO_SIT (10) | 60 |
| LIE_TO_STAND (12) | 58 |
| STAND_TO_LIE (11) | 58 |

Sit-to-stand has roughly **1/3 the instances** of the main static/walking activities, plus the transitions themselves are very short (~2.6s vs. much longer durations for static activities). This means both **instance count** and **total sample coverage** are heavily skewed against transitions.

Summary:

**Sit-to-stand is activity ID 8, not 5** as the checklist assumed. Important to hardcode 8 in your pipeline.

▼ **Explanation**

**The class imbalance is severe but at two levels:**

First, at the instance level: 62 sit-to-stand segments vs. 120+ for each static activity. That's manageable. But second, and more importantly, at the *sample/window* level: each sit-to-stand lasts ~2.6 seconds (~130 samples), while sitting/standing/laying segments can last tens of seconds each. When you window the data at 2.56s with 50% overlap, sit-to-stand will produce maybe 1–3 windows per instance (since the transition is roughly one window long), while a 30-second sitting segment produces ~22 windows. So you're looking at maybe ~100–180 sit-to-stand windows vs. thousands of "everything else" windows — probably a 5–10% positive class rate.

**What this means practically:**

The `class_weight='balanced'` setting in scikit-learn will handle this automatically by upweighting the minority class during training. You'll report precision, recall, and F1 for the sit-to-stand class specifically (not accuracy, which would be misleading — a model that always says "not sit-to-stand" would get 90%+ accuracy). The confusion matrix will be important for seeing what sit-to-stand gets confused with — expect stand-to-sit (activity 7) to be the main culprit since the signals are similar but reversed.

**Also notable:** each subject only has 2–3 sit-to-stand instances. That means in each LOSO-CV fold, the held-out subject will have very few positive test examples. This makes per-fold metrics noisy — some folds might show 100% recall, others 0%, depending on whether those 2–3 events happen to be easy or hard. The *mean across all 30 folds* will be the reliable number. Worth mentioning this in the writeup.

▼ **Phase 2: Feature Extraction Pipeline (Days 3-5)**

## Windowing

○ Write a function that takes a raw signal and cuts it into 2.56-scond with 50% overlap
  - At 50Hz, each window = 128 samples

- 50% overlap means each window starts 64 samples after the previous one
- Assign a label to each window based on which activity occupies the majority of that window
- **Drop ambiguous windows:** If a window straddles two different activity/transition boundaries (e.g., less than 80% of the window belongs to one label), drop it from the dataset. Transitions are exactly where labels are noisiest, and training on ambiguous labels hurts the model on the exact class you care about.
- Remap remaining labels to binary: sit-to-stand = 1, everything else = 0
- Track which subject each window belongs to (needed for LOSO-CV)

## Feature Computation

- First, compute **acceleration magnitude** per sample: `accel_mag = sqrt(ax² + ay² + az²)`. This is orientation-independent — "accel Z" is not guaranteed to be vertical because phone orientation varies in the waistband. Magnitude works regardless of how the phone is positioned.
- Similarly, compute **gyroscope magnitude**: `gyro_mag = sqrt(gx² + gy² + gz²)`
- For each window, compute 6 features for each of **8 channels** (accel x/y/z, gyro x/y/z, accel magnitude, gyro magnitude):

| Feature | Code | What it captures |
|---|---|---|
| Mean | `np.mean(window)` | Average signal level |
| Std | `np.std(window)` | Signal variability |
| Min | `np.min(window)` | Lowest value |
| Max | `np.max(window)` | Highest value |
| Range | `np.max(window) - np.min(window)` | Total swing |
| Energy | `np.mean(window**2)` | Movement intensity |

- This gives you **48 features per window** (6 features × 8 channels)
- Store as a DataFrame: each row = one window, columns = 48 features + label + subject_id
- Sanity check: print the shape — you should have thousands of rows, 50 columns

○ Sanity check: print label distribution — confirm sit-to-stand (label=1) is a small minority

## Verify Your Pipeline

○ Pick one known sit-to-stand segment from the labels file
○ Look at the raw signal for that segment — you should see a clear spike in acceleration magnitude
○ Confirm the windows covering that segment got labeled as 1
○ Pick one known sitting/standing segment — confirm those windows got labeled as 0
○ **This step catches bugs early. Do not skip it.**

▼ **Phase 2.1.1 - Windowing 1a**

   ▼ **Results**

### Script Structure

| Function | Purpose |
|----------|---------|
| `load_labels()` | Parses `labels.txt` → DataFrame with 5 int columns |
| `load_signals()` | Discovers all `acc_`/`gyro_` pairs, loads each with `np.loadtxt`, validates matching row counts, concatenates into `(T, 6)` arrays |
| `print_summary()` | Prints experiment/user/sample counts, per-activity segment counts, sit-to-stand duration stats, and per-subject breakdown |

The two return values are importable for downstream use:

```
from SRC.load_data import load_signals, load_labels
signals = load_signals()   # dict[(exp, user)] → ndarray (T, 6)
labels  = load_labels()    # DataFrame
```

### Key Findings

- **61 experiment-user pairs**, 30 subjects, 1.12M total samples (~6.2 hours)

- **No row mismatches** between any acc/gyro pair — all clean
- **62 sit-to-stand segments** (activity ID 8) — most subjects have exactly 2; users 8 and 22 have 3
- **Average transition duration: 129.5 samples = 2.59 seconds** (range 1.48–3.66s)
- **Severe class imbalance**: transitions like sit-to-stand have ~62 segments vs. 120–186 for static/walking activities, and the transitions are much shorter in duration

▼ **Phase 2.1.2 - Windowing 2a**

 ▼ **Results**

It runs perfectly. The output reveals something important — several users (7, 9, 10, 12, 27) have **0 positive windows** despite having sit-to-stand segments. That's because their sit-to-stand transitions were too short to fill 80% of a 128-sample window. Let me verify this is correct and not a bug:

The logic is correct. Users 7, 10, and 27 genuinely have sit-to-stand segments too short (< 103 samples) to ever pass the 80% purity threshold in a 128-sample window. Users 9 and 12 have segments long enough in theory but the window stride alignment causes them to straddle boundaries. This is expected behavior — the script is working correctly.

Here's the summary of what `src/windowing.py` does and the results:

## Script API

```python
from src.windowing import create_windows

X, y, users, n_dropped = create_windows(signals, labels_df)
# X:          (15321, 128, 6) — windowed sensor data
# y:          (15321,)        — binary labels (1=sit-to-stand,
0=other)
# users:      (15321,)        — subject IDs for LOSO-CV
# n_dropped: 2132             — ambiguous windows discarded
```

## Key Results

| Metric | Value |
| --- | --- |
| Windows kept | 15,321 |
| Sit-to-stand (label=1) | 51 |
| Other (label=0) | 15,270 |
| Positive class % | **0.33%** |
| Dropped (ambiguous) | 2,132 |

The class imbalance is extreme — **300:1 negative-to-positive ratio**. Five users (7, 9, 10, 12, 27) end up with zero positive windows because their sit-to-stand transitions were shorter than the ~103 samples needed to clear the 80% purity threshold within a 128-sample window. This is something to keep in mind when designing the LOSO-CV splits and choosing evaluation metrics (precision/recall/F1 over accuracy).

▼ Explanation

Correct. Out of 15,321 windows (each 2.56 seconds long), only 51 are labeled as sit-to-stand. The rest are walking, sitting, standing, laying, other transitions, or unlabeled gaps — all grouped as "everything else."

51 is low. It's a direct consequence of sit-to-stand being a very brief event (~2.6 seconds) that only happens 62 times in the whole dataset, combined with the 80% purity threshold dropping windows that straddle boundaries.

Is something about this concerning you? We can proceed and see how the models do, or if you're worried about having too few positive examples, we could discuss options like lowering the purity threshold to recover more windows.

The purity threshold is the 80% rule from the checklist. When a window straddles two different activities — say the first half is "sitting" and the second half is "sit-to-stand" — you have to decide what label to give it. The 80% rule says: if less than 80% of the samples in that window belong to a single activity, drop the window entirely.

**Why it exists:** A window that's 50% sitting and 50% sit-to-stand is ambiguous. If you label it as sit-to-stand, you're training the model on a signal that's half sitting. If you label it as "other," you're teaching the model that a signal containing real sit-to-stand movement is negative. Either way you're feeding the model noisy data. Dropping it avoids both problems.

**The cost:** Since sit-to-stand transitions are only ~130 samples and windows are 128 samples, many sit-to-stand windows end up straddling the boundary and getting dropped. That's why you went from 62 segments down to 51 windows — and why 5 users lost all their positive windows entirely.

The 2,132 dropped windows in your output are all windows (across all activities) that failed the 80% purity check. Most of those are at boundaries between any two activities, not just sit-to-stand.

▼   **Phase 2.1.2 - Windowing 2b**

    ▼   **Results**

## All 62 Sit-to-Stand Segments

The full table is printed above. The  ∗  flag marks segments too short to ever produce a pure window.

## Duration Statistics

| Stat | Samples | Seconds |
|------|---------|---------|
| Min | 74 | 1.48 |
| Max | 183 | 3.66 |
| Mean | 129.5 | 2.59 |
| Median | 129.5 | 2.59 |
| Std | 27.7 | 0.55 |

## Purity Threshold Impact

- **10 / 62 segments (16.1%)** are shorter than 103 samples — these can **never** produce a window that passes 80% purity in a 128-sample window
- The affected users: **7** (both segments: 93, 74), **8** (one of three: 74), **9** (one of two: 97), **10** (both: 89, 78), **22** (one of three: 92), **26** (one of two: 102), **27** (both: 98, 101)
- Users 7, 10, and 27 lose **all** their sit-to-stand segments — they will have zero positive windows in LOSO-CV, meaning those folds can only evaluate specificity, never recall

This is worth keeping in mind when evaluating LOSO-CV results: 3 of the 30 folds are structurally unable to detect any true positives.

Here's the comparison:

## 50% vs 80% Purity Threshold

| Metric | 80% threshold | 50% threshold | Change |
|---|---|---|---|
| **Windows kept** | 15,321 | 17,453 | +2,132 |
| **Sit-to-stand (1)** | 51 | 126 | **+75 (2.5×)** |
| **Other (0)** | 15,270 | 17,327 | +2,057 |
| **Positive class %** | 0.33% | 0.72% | **~2× richer** |
| **Windows dropped** | 2,132 | 0 | **Zero drops** |

Key improvements:

- **Zero windows dropped** — with a 50% threshold, the majority vote always has a winner (ties at exactly 50/50 between two activities are essentially impossible given the label structure).

- **All 30 users now have positive windows** — the previously-starved users 7, 9, 10, 12, and 27 now each have 3–4 sit-to-stand windows. No user is left out.

- The positive class roughly doubled from 0.33% to 0.72% — still heavily imbalanced, but every LOSO-CV fold now has at least some true positives to evaluate recall against.

▼  **Explanation**

**One deviation from checklist to note in your writeup:** You lowered the purity threshold from 80% to 50%, supported by data showing 16% of sit-to-stand segments were too short to ever pass 80%. This recovered 75 additional positive windows (51 → 126) and ensured all 30 subjects have positive examples for LOSO-CV.

▼  **Phase 2.2 - Feature Computation**

▼  **Result**

Confirmed — the SITTING segment (activity 4, samples 1393–2194) is immediately followed by the SIT_TO_STAND segment (activity 8, samples 2195–2359). The 2 red bars in the bottom panel's right edge are sit-to-stand windows that fall within the extended context region. This is exactly correct.

## Summary

`src/features.py` is complete and produces:

### Feature Matrix — `Data/uci_hapt/features.csv`

| Property | Value |
|---|---|
| **Shape** | 17,453 rows × 50 columns |
| **Feature columns** | 48 (8 channels × 6 stats) |
| **Metadata columns** | `label`, `subject_id` |
| **label=0 (other)** | 17,327 (99.28%) |
| **label=1 (sit-to-stand)** | 126 (0.72%) |

### 8 Channels

`ax, ay, az, gx, gy, gz` (raw) + `accel_mag, gyro_mag` (computed via L2 norm)

### 6 Features per Channel

`mean, std, min, max, range, energy` → named like `ax_mean`, `gyro_mag_energy`, etc.

### Sanity Check Plot — `Results/sanity_check.png`

- **Top panel**: Exp 1, User 1 — the first sit-to-stand segment. The characteristic acceleration spike is clearly visible. **2 windows correctly labelled as positive** (red bars) overlap the orange ground-truth region.
- **Bottom panel**: Exp 1, User 1 — a sitting segment. Flat, low-amplitude signal. **All windows within the blue region correctly labelled as negative** (green bars). The 2 red bars at the far right are the adjacent sit-to-stand windows from the context margin — confirming the labelling pipeline is consistent end-to-end.

▼ **Explanation**

Let me explain what the sanity check plot is showing and why it matters.

**What the plot does:** It picks two segments from the same subject and shows the raw acceleration magnitude signal with colored overlays indicating the ground truth labels and your pipeline's window labels.

**Top panel (sit-to-stand segment):** The signal shows a clear spike — that's the person pushing up from the chair. The acceleration jumps because the body is accelerating upward. The orange shaded region is the ground truth: the labels file says "sit-to-stand happens here." The red bars are windows that your pipeline labeled as 1 (sit-to-stand). Two windows overlap that orange region and both got labeled correctly as positive. This confirms your pipeline is assigning the right label to the right windows.

**Bottom panel (sitting segment):** The signal is flat and low-amplitude — the person is just sitting still, not much movement. The blue shaded region is the ground truth for sitting. The green bars are windows your pipeline labeled as 0 (not sit-to-stand). All correct. The two red bars at the far right edge aren't a bug — they're from the sit-to-stand segment that happens immediately *after* the sitting segment. The plot shows a little extra context beyond the sitting region, and those red windows fall in that adjacent sit-to-stand zone. This actually confirms the pipeline is working: the sitting windows are green, and the moment the sit-to-stand starts, the labels flip to red.

**Why this matters:** This is a visual proof that your entire pipeline — loading data, aligning labels to samples, windowing, majority-vote labeling — is correct. If something was off (wrong sample alignment, off-by-one error, flipped labels), it would show up immediately in this plot. The checklist says "do not skip this step" because bugs in the feature pipeline silently corrupt everything downstream.

**Phase 2 is now complete.** You have a clean 17,453 × 50 feature matrix ready for modeling. Want to move to Phase 3 — building the three models and running a quick sanity test?

▼ **Feature Matrix**

12

Here's an interactive view of what your feature matrix looks like. The key things to notice:

**The orange-highlighted row (row 4,831)** is a sit-to-stand window — label = 1. Compare its numbers to the rows above it (sitting, label = 0). The sit-to-stand window has much higher acceleration magnitude max (2.31g vs ~1.05g) and much higher gyroscope activity. That's the signal the model will learn to recognize.

**The ··· columns** represent the ~36 features we can't fit on screen. The full matrix has 48 feature columns plus label and subject_id — 50 columns total.

The model's job is simple: look at a row of 48 numbers and decide if it's a 1 or a 0.

▼ **Phase 3: Build the Three Models (Days 5-7)**

### Model 1: Threshold Baseline (Non-ML)

○ This is your "no learning" approach

○ Logic: use **acceleration magnitude** features (not a single axis — orientation-independent)

○ If `max_accel_mag > threshold` and `range_accel_mag > threshold`, predict sit-to-stand

○ Tune thresholds by eyeballing the data or using a simple grid search on training data

○ Note: this baseline should be simple and dumb on purpose — it represents what you'd build without ML

### Model 2: Logistic Regression

○ `from sklearn.linear_model import LogisticRegression`
○ Use `class_weight='balanced'` to handle class imbalance
○ Train on your 48-feature matrix with binary labels
○ This is your "simple linear ML" approach

### Model 3: Random Forest

○ `from sklearn.ensemble import RandomForestClassifier`
○ Use `class_weight='balanced'` to handle class imbalance
○ Start with default hyperparameters (100 trees)
○ This is your "ensemble ML" approach

### Quick Test (before full LOSO-CV)

○ Do a quick sanity check: train on subjects 1-25, test on subjects 26-30

- ○ Print precision, recall, F1 for class 1 (sit-to-stand)
- ○ If all three models give 0% recall, something is wrong with your pipeline — debug before proceeding
- ○ If numbers look reasonable (>50% recall), proceed to full LOSO-CV

▼ **Results**

All three models are well above 0% recall — the pipeline is working. Here's the summary:

## Sanity Test Results (Train: users 1–25, Test: users 26–30)

|  | Train pos | Test pos |
| --- | --- | --- |
| **Split** | 106 / 14,393 | 20 / 3,060 |

| Model | Precision | Recall | F1 | TP | FP | FN |
| --- | --- | --- | --- | --- | --- | --- |
| **Threshold Baseline** | 0.010 | **0.950** | 0.020 | 19 | 1,821 | 1 |
| **Logistic Regression** | 0.179 | **0.950** | 0.302 | 19 | 87 | 1 |
| **Random Forest** | **1.000** | 0.650 | **0.788** | 13 | 0 | 7 |

**What this tells us:**

- ○ **Threshold Baseline** finds nearly every sit-to-stand window (95% recall) but fires on 1,821 false positives — it's a blunt instrument with just 2 features and no way to suppress false alarms.
- ○ **Logistic Regression** also hits 95% recall but with far fewer false positives (87), giving a much more useful F1 of 0.30.
- ○ **Random Forest** has perfect precision (zero false positives) at the cost of missing 7 of 20 positives (65% recall), yielding the best F1 at 0.79.

All three are clearly learning something meaningful — ready to proceed to full LOSO-CV.

▼ **Explanation**

Let me walk through each model and what the results mean.

---

## Model 1: Threshold Baseline

**How it works:** This is the simplest possible approach — no machine learning at all. It looks at just 2 numbers from each window:

- `accel_mag_max` : the highest acceleration magnitude in that 2.56 seconds
- `accel_mag_range` : the difference between highest and lowest acceleration magnitude

If both numbers are above a threshold, it says "sit-to-stand." That's it. The thresholds were found by trying many combinations on the training data and picking what gave the best F1.

**Why it exists:** It represents what you could build with zero ML knowledge — just "sit-to-stand involves a big acceleration spike, so look for big spikes." It's the baseline the ML models need to beat to justify their complexity.

**Results: Precision 0.01, Recall 0.95, F1 0.02**

It found 19 of the 20 real sit-to-stands (great recall). But it also flagged 1,821 other windows as sit-to-stand (terrible precision). Why? Because lots of activities produce big acceleration spikes — walking upstairs, walking downstairs, standing up from lying down. With only 2 features, the model can't tell these apart. It's like a smoke detector that goes off every time you cook — it catches real fires, but also gives constant false alarms.

**The 0.01 precision means:** out of every 100 windows it labels as sit-to-stand, only 1 actually is. Useless in practice.

---

## Model 2: Logistic Regression

**How it works:** This is the simplest ML model. It takes all 48 features, multiplies each by a learned weight, adds them up, and if the total is above zero, predicts sit-to-stand. Think of it as a formula:

```
score = (w1 × ax_mean) + (w2 × ax_std) + ... + (w48 ×
gyro_mag_energy) + bias
if score > 0: predict sit-to-stand
```

During training, it finds the weights that best separate sit-to-stand from everything else. `class_weight='balanced'` tells it to care about getting sit-to-stand right even though there are very few examples. `StandardScaler` puts all features on the same scale so no single feature dominates just because it has bigger numbers.

**The key limitation:** the decision boundary is *linear* — it can only draw a straight line (or hyperplane in 48 dimensions) between the two classes. It can't learn rules like "high acceleration AND high gyroscope means sit-to-stand, BUT high acceleration AND low gyroscope means walking."

**Results: Precision 0.179, Recall 0.95, F1 0.30**

Same 95% recall as the threshold baseline — it still finds 19 of 20 positives. But false positives dropped from 1,821 to 87. That's a massive improvement. By using all 48 features instead of just 2, it can rule out activities that have high acceleration but different gyroscope patterns, different axis distributions, etc.

**The 0.179 precision means:** out of roughly every 6 windows it labels as sit-to-stand, 1 actually is. Better than 1-in-100, but still noisy. The model is still biased toward saying "yes" because `class_weight='balanced'` heavily upweights the rare positive class.

---

## Model 3: Random Forest

**How it works:** It builds 100 decision trees, each one a series of if/else rules:

```
if accel_mag_max > 1.8:
    if gyro_mag_std > 0.5:
        if ax_range > 0.7:
            predict sit-to-stand
        else:
            predict other
    else:
        predict other
else:
    predict other
```

Each tree is trained on a random subset of the data and a random subset of the features. The final prediction is a majority vote across all 100 trees. This makes it robust — individual trees might overfit, but the ensemble averages out the noise.

**Results: Precision 1.00, Recall 0.65, F1 0.79**

The opposite pattern from the other two models. When it says sit-to-stand, it's always right (zero false positives). But it misses 7 of 20 real sit-to-stands. It's being conservative — only predicting sit-to-stand when it's very confident.

**The 1.00 precision means:** every window it labels as sit-to-stand actually is one. But 0.65 recall means it's silent on 35% of real sit-to-stands. In an app context, this means it would undercount reps.

## What the comparison tells you

|  | Threshold | Logistic Regression | Random Forest |
|---|---|---|---|
| **Approach** | 2 hand-picked rules | Linear boundary, 48 features | 100 nonlinear trees, 48 features |
| **Finds sit-to-stands?** | Almost all | Almost all | Most |
| **False alarms?** | Massive | Moderate | Zero |
| **Best F1** | 0.02 | 0.30 | 0.79 |

The progression from left to right shows that the problem *does* benefit from ML and *does* benefit from nonlinear feature interactions. The threshold baseline proves the signal exists. Logistic regression shows that more features help. Random Forest shows that nonlinear decision rules dramatically improve precision.

**For your app**, you want both high precision (don't count fake reps) and high recall (don't miss real reps). Random Forest is winning on F1 because it balances these well. In Phase 4 (LOSO-CV across all 30 subjects), you'll get the real numbers — this quick test was just to confirm nothing is broken.

Does this make sense? Any part you want me to dig into further?

**The key advantage over logistic regression:** it can learn *nonlinear* interactions. "High acceleration AND high gyroscope" can mean something different from "high acceleration AND low gyroscoe." Logistic regression can't do this natively.

▼ **Phase 4: LOSO-CV Evaluation — Experiment 1 (Days 7-9)**

## Implement LOSO-CV

○ Loop through subjects 1 to 30
○ In each fold:
- Train set = all windows NOT from this subject
- Test set = all windows from this subject
- Train each of the 3 models on train set
- Predict on test set
- Record precision, recall, F1 for sit-to-stand class

- Record PR-AUC (precision-recall area under curve) for ML models — `from sklearn.metrics import average_precision_score`. Use predicted probabilities, not hard labels. PR-AUC is more informative than F1 for rare classes.
- ○ Store results for each fold in a list/DataFrame

## Compute and Report Window-Level Results

- ○ Calculate mean ± std across 30 folds for precision, recall, F1, PR-AUC — for each of the 3 models
- ○ Create a comparison table:

| Model | Precision | Recall | F1 | PR-AUC |
|---|---|---|---|---|
| Threshold baseline | X ± Y | X ± Y | X ± Y | — |
| Logistic regression | X ± Y | X ± Y | X ± Y | X ± Y |
| Random Forest | X ± Y | X ± Y | X ± Y | X ± Y |

- ○ Generate a confusion matrix (aggregated across all folds) for the best model
- ○ Check: what are sit-to-stand windows most commonly confused with? (probably stand-to-sit)

## Rep-Level Evaluation (Event Detection)

- ○ For each held-out subject in each fold, apply post-processing to turn window predictions into events:
  1. **Smooth predictions:** Apply a moving average (0.5–1.0s window) to predicted probabilities, then threshold
  2. **Minimum event duration:** Discard any positive cluster shorter than 0.5s (too short to be a real sit-to-stand)
  3. **Minimum gap between events:** Merge positive clusters separated by less than 1.0s (likely the same transition)
  4. **Count remaining clusters** = predicted rep count
- ○ Define event-level matching: a predicted event is a true positive if it falls within ±0.5s of a labeled sit-to-stand transition time. This gives you event-level precision and recall.
- ○ Compare predicted rep count against ground truth rep count per subject

○ Report mean absolute error between predicted and actual rep counts

○ Report event-level precision and recall alongside window-level metrics

○ This is the metric that connects directly to the app — and is more meaningful for your story

### Feature Importance (RF only)

○ After full LOSO-CV, train a final RF on all 30 subjects

○ Print `model.feature_importances_`

○ Plot top 10 features

○ Check: do the top features align with the literature? (expect acceleration magnitude and vertical features to rank high, per Millor and Galán-Mercant)

○ Discuss in writeup

▼ **Phase 4.1. - LOSO-CV + Window-Level Results**

  ▼ **Results**

Everything ran successfully. Here's the full breakdown:

## LOSO-CV Results

### Aggregate Summary (30 folds)

| Model | Precision | Recall | F1 | PR-AUC |
|---|---|---|---|---|
| **Threshold Baseline** | 0.011 ± 0.002 | 0.933 ± 0.135 | 0.022 ± 0.004 | — |
| **Logistic Regression** | 0.186 ± 0.070 | 0.933 ± 0.117 | 0.305 ± 0.096 | 0.675 ± 0.148 |
| **Random Forest** | 0.778 ± 0.377 | 0.454 ± 0.283 | **0.554 ± 0.301** | **0.742 ± 0.212** |

All 30 folds have at least 1 positive test example (thanks to the 50% purity threshold), so `Recall (all)` and `Recall (pos>0)` are identical.

### Best Model: Random Forest (mean F1 = 0.554)

**Pooled confusion matrix:**

- **TP = 58** — correctly detected sit-to-stand windows

- **FP = 5** — false alarms (only 5 across all 17,453 windows!)

- **FN = 68** — missed sit-to-stand windows

- **TN = 17,322**

## Hardest Subjects

Subjects **9, 10, 17, 23, 25** got 0 TP from the Random Forest — complete misses. These are subjects with short sit-to-stand segments that produce noisy boundary windows. Subject **4** also struggled (1/5 recall = 20%).

Best folds: subjects **24** and **27** hit perfect 1.000 F1.

## False-Positive Breakdown

**100% of Random Forest false positives were SITTING windows.** All 5 FP windows had true activity = SITTING (activity 4). This makes physical sense — sitting is the activity immediately before a sit-to-stand transition, and the boundary windows between them share similar sensor patterns.

## Model Comparison Interpretation

- **Threshold Baseline**: Catches 93% of sit-to-stand but fires on ~60% of all windows as positives — completely unusable precision (1.1%).
- **Logistic Regression**: Same 93% recall but precision jumps to 19% — 17× better than threshold. Still lots of false alarms.
- **Random Forest**: Flips the tradeoff — 78% precision (almost no false alarms) at the cost of catching only 45% of sit-to-stand events. Best F1 and best PR-AUC.

## Outputs

- `Results/loso_cv_results.csv` — all 90 rows (30 subjects × 3 models)
- `Results/confusion_matrix.png` — pooled confusion matrix for Random Forest

▼ **Explanation**

Great results. Let me explain what these mean before you run Parts 2 and 3.

**The big picture:** Random Forest wins on F1 (0.554) and PR-AUC (0.742), confirming that the problem benefits from nonlinear feature interactions. But the results also reveal a clear tension between precision and recall that none of the models fully resolve.

**Random Forest: 78% precision, 45% recall.** It's very conservative — when it says sit-to-stand, it's almost always right (only 5 false positives across the entire dataset). But it misses more than half of the actual sit-to-stand windows. The high variance across folds (F1 std of 0.301) tells you performance is inconsistent — some subjects are easy (perfect F1), others are completely missed.

**Why 5 subjects get 0 TP:** These subjects likely have sit-to-stand transitions that look different from the majority — shorter duration, less forceful movement, or different body mechanics. The model trained on the other 29 subjects never learned to recognize their pattern. This is actually a preview of the generalization challenge you'll face with elderly populations, who move differently from the 19–48 year olds in UCI HAPT.

**Why all 5 false positives are SITTING:** This is a great finding for your writeup. The boundary between sitting and sit-to-stand is where the signal transitions from low activity to high activity. A window that captures the very end of sitting (person starting to lean forward before standing) looks similar to the beginning of sit-to-stand. This is the exact boundary that the 50% purity threshold makes noisier — a good tradeoff to discuss.

**The PR-AUC numbers matter:** RF gets 0.742 vs Logistic Regression's 0.675. PR-AUC evaluates across all possible probability thresholds, not just the default 0.5 cutoff. The RF's higher PR-AUC means it has better overall ranking ability — it assigns higher probabilities to true sit-to-stand windows than to non-sit-to-stand windows. This suggests that with threshold tuning, you could potentially trade some precision for better recall.

**One important thing for the writeup:** the 45% recall at window level might look concerning, but window-level recall isn't the same as rep-level recall. The model might miss some windows within a sit-to-stand event but still detect enough consecutive windows to count the rep. That's exactly what Part 2 (event detection) will reveal — and it's likely to tell a more optimistic story.

▼ **Phase 4.2 - Rep-Level Evaluation (Event Detection)**

   ▼ **Results 1**

## Event-Level Summary

| Model | Rep MAE | Event Precision | Event Recall | Event F1 |
|---|---|---|---|---|
| **Threshold Baseline** | 16.27 | 0.091 | **0.806** | 0.163 |
| **Logistic Regression** | 3.23 | 0.151 | 0.387 | **0.217** |
| **Random Forest** | **1.63** | **0.538** | 0.113 | 0.187 |

## Pooled Counts (all 30 folds)

| Model | TP | FP | FN | GT Total |
|---|---|---|---|---|
| Threshold Baseline | 50 | 500 | 12 | 62 |
| Logistic Regression | 24 | 135 | 38 | 62 |
| Random Forest | 7 | 6 | 55 | 62 |

## What's happening

**Threshold Baseline** predicts ~16–30 events per subject (GT is 2–3). It finds most real events (81% recall) but drowns them in false alarms — averaging 16.7 extra events per subject.

**Logistic Regression** is the best balanced model at the event level (highest event F1 = 0.217). It over-counts by ~3 reps on average and finds 39% of real events.

**Random Forest** has the best rep-count accuracy (MAE = 1.63) and the fewest false alarms (only 6 total across all subjects). But the min–event–duration filter of 2 windows is aggressive for its already-sparse predictions — it only detects 7 of 62 events. The RF tends to correctly predict just 1 or 2 isolated positive windows per transition, which often get discarded as noise.

## Post-processing parameters used

- Smoothing: moving average k=3 on probabilities (LR/RF only), threshold at 0.5
- Minimum event duration: 2 consecutive positive windows
- Maximum merge gap: 1 window
- Match tolerance: ±1.0 seconds

## Output

- `Results/event_detection_results.csv` — 90 rows (30 subjects × 3 models), 11 columns

▼ **Explanation 1**

These results reveal an important problem. The post-processing is too aggressive for Random Forest's prediction pattern. RF often correctly identifies a sit-to-stand but only predicts 1 positive window for it (because the transition is short and RF is conservative). Then the minimum event duration filter of 2 consecutive windows throws it away as noise. You're discarding correct predictions.

**The fix is to relax the post-processing for RF specifically.**

The key insight here: RF's strength is precision (it rarely false-alarms), so it doesn't need aggressive filtering to remove false positives. The strict post-processing was designed for noisy models like the threshold baseline. RF needs lighter post-processing that preserves its sparse but accurate predictions.

After this runs, share all three tables and we'll pick the best configuration. This is also a good discussion point for your writeup — how post-processing interacts differently with different model types.

▼ **Results 2**

The results are dramatic. Here's the full comparison:

## Side-by-Side: Three Post-Processing Configs

### Threshold Baseline

| Config | MAE | Precision | Recall | F1 | TP | FP | FN |
|--------|-------|-----------|--------|-------|----|-----|----|
| Strict | 16.27 | 0.091 | 0.806 | 0.163 | 50 | 500 | 12 |
| Relaxed | 14.73 | 0.101 | 0.823 | 0.180 | 51 | 453 | 11 |
| Minimal | 14.73 | 0.101 | 0.823 | 0.180 | 51 | 453 | 11 |

Post-processing barely helps — this model is so noisy that even aggressive filtering can't clean it up.

### Logistic Regression

| Config | MAE | Precision | Recall | F1 | TP | FP | FN |
|--------|-------|-----------|--------|-------|----|-----|----|
| Strict | 3.23 | 0.151 | 0.387 | 0.217 | 24 | 135 | 38 |
| Relaxed | 4.03 | 0.142 | 0.419 | 0.212 | 26 | 157 | 36 |
| **Minimal** | **11.07** | **0.155** | **0.984** | **0.268** | **61** | **333** | **1** |

Removing smoothing unlocks near-perfect recall (61/62 events found!) but floods with false positives. Smoothing acts as a crucial noise filter for LR.

## Random Forest

| Config | MAE | Precision | Recall | F1 | TP | FP | FN |
|--------|-----|-----------|--------|-----|-----|-----|-----|
| Strict | 1.63 | 0.538 | 0.113 | 0.187 | 7 | 6 | 55 |
| Relaxed | 1.50 | 0.588 | 0.161 | 0.253 | 10 | 7 | 52 |
| **Minimal** | **0.60** | **0.909** | **0.645** | **0.755** | **40** | **4** | **22** |

## Winner: Random Forest + Minimal post-processing

**Event F1 = 0.755** — a massive jump from 0.187 under strict settings. The numbers tell the story:

- **91% precision** — only 4 false alarms across all 30 subjects
- **65% recall** — finds 40 of 62 real sit-to-stand events
- **MAE = 0.60 reps** — the average rep-count error is less than 1
- **20 of 30 subjects** got a perfect rep count (|error| = 0)
- Subjects 9, 10, 17, 23, 25 remain at 0 detections — these have the shortest STS segments

This confirms the hypothesis exactly: RF's sparse but accurate single-window predictions were being destroyed by the min-duration=2 filter. With minimal post-processing, those lone correct predictions survive and become correctly matched events. Saved to `Results/event_detection_results_tuned.csv`.

▼ **Explanation 2**

Let me walk through each model and configuration so you can see what happened.

---

## What the three post-processing configs do

**Strict** (original): Smooths the predictions, requires at least 2 consecutive positive windows to count as an event, merges clusters that are 1 window apart. This is aggressive — it throws away isolated single positive windows as noise.

**Relaxed**: Same smoothing, but a single positive window can count as an event, and clusters 2 windows apart get merged. Less aggressive.

**Minimal**: No smoothing at all, single positive windows count as events, clusters 2 windows apart get merged. Basically trusts whatever the model predicts.

---

## Threshold Baseline — post-processing doesn't help

The threshold baseline produces a flood of positives — it flags roughly 60% of all windows as sit-to-stand. No amount of post-processing can clean that up. Strict filtering removes some false alarms but also loses real events. Relaxed/minimal settings barely change anything because the false positives are so dense and continuous that they form large clusters that survive any filter. The model itself is too noisy for post-processing to rescue.

---

## Logistic Regression — smoothing is critical

The interesting row is **Minimal**: recall jumps to 98% (finds 61 of 62 events) but false positives triple to 333. Why? Without smoothing, every brief spike in predicted probability above 0.5 becomes an event. LR produces lots of these brief spikes — windows near activity transitions where the probability fluctuates around 0.5. Smoothing (averaging 3 consecutive predictions) dampens these fluctuations, which is why the **Strict** config works better for LR. It needs that noise filter.

The tradeoff: strict LR finds fewer events (24/62) but with fewer false alarms (135). It's mediocre either way — too many false positives to be useful, too many misses to be reliable.

---

## Random Forest — minimal wins dramatically

This is the key finding. Under **Strict** settings, RF only detects 7 of 62 events despite having good window-level predictions. Here's exactly what was happening:

1.  A sit-to-stand transition occurs
2.  RF correctly predicts 1 positive window for it (because the transition is short — ~2.6 seconds ≈ 1–2 windows)
3.  The smoothing step averages that single positive prediction with the surrounding negative predictions, pulling the probability below 0.5
4.  Even if it survives smoothing, the min_duration=2 filter discards it because it's only 1 window long
5.  A correct prediction gets thrown away

Under **Minimal** settings:

1. Same sit-to-stand transition occurs
2. RF correctly predicts 1 positive window
3. No smoothing — the prediction survives as-is
4. Min_duration=1 — a single positive window counts as an event
5. The correct prediction becomes a detected event

The reason this works for RF but not for LR or the threshold baseline: **RF rarely makes false positive predictions** (only 5 across the entire LOSO-CV at window level). So there's no noise to filter out. When RF says "sit-to-stand," it's almost always right. The strict post-processing was solving a problem RF doesn't have (false alarms) while creating a problem it does have (discarding sparse correct predictions).

---

## The final numbers for RF + Minimal

Think about what these mean for the app:

**MAE = 0.60 reps**: If someone does 2 sit-to-stands, the model says 2. If someone does 3, the model usually says 3, occasionally says 2. On average, it's off by less than 1 rep. 20 of 30 subjects get an exact count.

**91% event precision**: When the app draws a marker saying "rep detected here," 91% of the time there really was a sit-to-stand at that moment. Only 4 phantom reps across all 30 subjects.

**65% event recall**: The app catches about 2 out of every 3 real sit-to-stands. The 35% it misses are concentrated in 5 specific subjects with unusually short or gentle transitions.

**The 5 missed subjects**: These aren't random failures — they're subjects whose movement pattern is too different from the majority. In your writeup, this connects to the generalization concern: if the model struggles with unusual movers within the training population (young adults), it may struggle more with elderly populations who move differently. That's exactly what Experiment 2 (external validation) will test.

---

Does this all make sense? Ready for Part 3 (feature importance)?

▼ **Phase 4.3 - Feature Importance**

Everything looks clean. Here's the summary:

## Feature Importance Results

### Top 10 Features (by Gini importance)

| Rank | Feature | Importance | Group |
|------|---------|------------|-------|
| 1 | **gx_energy** | 0.0937 | Gyro axes |
| 2 | **accel_mag_std** | 0.0755 | Accel magnitude |
| 3 | gx_range | 0.0696 | Gyro axes |
| 4 | ax_std | 0.0680 | Accel axes |
| 5 | gx_std | 0.0676 | Gyro axes |
| 6 | gz_mean | 0.0524 | Gyro axes |
| 7 | ax_range | 0.0478 | Accel axes |
| 8 | accel_mag_max | 0.0448 | Accel magnitude |
| 9 | gx_max | 0.0396 | Gyro axes |
| 10 | accel_mag_energy | 0.0375 | Accel magnitude |

### Channel Group Dominance (top-10)

| Group | Count | % of importance |
|-------|-------|-----------------|
| **Gyro axes** | 5 features | **54.1%** |
| Accel magnitude | 3 features | 26.5% |
| Accel axes | 2 features | 19.4% |

### Key Findings

**Gyroscope x-axis dominates** — `gx_energy` is the single most important feature. This makes physical sense: sit-to-stand involves a forward rotation around the x-axis (pitch), generating a strong gyroscope signal that distinguishes it from static activities like sitting or standing.

**Literature alignment: Confirmed.** Three of four expected accel-magnitude features (`accel_mag_std`, `accel_mag_max`, `accel_mag_energy`) appear in the top 10. `accel_mag_range` is rank 13 — just outside. This matches Millor et al. (2014).

**Gyroscope is the surprise winner.** Five gyro features in the top 10 (54% of importance) confirms Galán-Mercant et al.'s finding that rotational dynamics provide critical complementary information. This also supports the value of Experiment 3 in your project plan (accel-only vs full feature ablation) — removing gyroscope features would eliminate the model's most discriminative signals.

▼

▼ Explanation

**The gyroscope dominance is your most interesting result from feature importance.** The model relies most heavily on rotational signals — especially around the x-axis (pitch). This makes biomechanical sense: sit-to-stand involves the trunk rotating forward and then upward, which produces a distinctive gyroscope signature that walking, sitting, or other activities don't share. Acceleration spikes alone are ambiguous (walking upstairs also produces high acceleration), but the combination of acceleration spike + specific rotational pattern is what makes sit-to-stand identifiable.

**This directly supports two things in your project:**

First, it validates the literature. Millor et al. found angular velocity was the strongest differentiator across frailty levels. Your model independently discovered the same thing — gyroscope features are the most informative. This is worth a sentence in your discussion.

Second, it foreshadows a problem for Experiment 3 (feature ablation). If gyroscope provides 54% of the model's discriminative power, the accel-only version will likely take a significant performance hit. And since the Marques dataset (elderly subjects) only has accelerometer data, your external validation on that dataset will be working without the model's best features.

**Phase 4 is now complete:**

- ✅ LOSO-CV with window-level results for all 3 models
- ✅ Confusion matrix — false positives are all SITTING windows
- ✅ Event-level evaluation with post-processing tuning
- ✅ Best config: RF + minimal post-processing (MAE 0.60, event F1 0.755)
- ✅ Feature importance — gyroscope x-axis dominates, aligns with literature

▼ Phase 5: Feature Ablation — Experiment 3 (Days 9-10)

## Accelerometer-Only vs. Full

○ Re-run LOSO-CV using only accelerometer features (24 features: 6 stats × 3 accel axes + accel magnitude)

○ Compare against the full 48-feature results

○ Create comparison table:

| Feature Set | Precision | Recall | F1 | PR-AUC |
|---|---|---|---|---|
| Accel only (24 features) | X ± Y | X ± Y | X ± Y | X ± Y |
| Accel + Gyro (48 features) | X ± Y | X ± Y | X ± Y | X ± Y |

○ This answers: does gyroscope add value?

○ Practically important because Marques only has accelerometer data

▼ **Results**

The results are very clear. Here's the full comparison:

## Window-Level: Accel-Only vs Full Features

| Feature Set | Model | Precision | Recall | F1 | PR-AUC |
|---|---|---|---|---|---|
| **Accel only (24)** | Threshold Baseline | 0.011 ± 0.002 | 0.933 ± 0.135 | 0.022 ± 0.004 | — |
| **Accel only (24)** | Logistic Regression | 0.077 ± 0.024 | 0.903 ± 0.150 | 0.141 ± 0.041 | 0.401 ± 0.156 |
| **Accel only (24)** | Random Forest | 0.306 ± 0.431 | 0.109 ± 0.161 | 0.152 ± 0.214 | 0.445 ± 0.234 |
| **Full (48)** | Threshold Baseline | 0.011 ± 0.002 | 0.933 ± 0.135 | 0.022 ± 0.004 | — |
| **Full (48)** | Logistic Regression | 0.186 ± 0.070 | 0.933 ± 0.117 | 0.305 ± 0.096 | 0.675 ± 0.148 |
| **Full (48)** | Random Forest | 0.778 ± 0.377 | 0.454 ± 0.283 | 0.554 ± 0.301 | 0.742 ± 0.212 |

**Δ (Full – Accel-only) for Random Forest: +0.472 precision, +0.345 recall, +0.402 F1, +0.296 PR-AUC**

## Event-Level: RF + Minimal Post-Processing

| Feature Set | Rep MAE | Event Precision | Event Recall | Event F1 | TP | FP | FN |
|---|---|---|---|---|---|---|---|
| Accel only (24) | 1.53 | 0.700 | 0.226 | 0.341 | 14 | 6 | 48 |
| **Full (48)** | **0.60** | **0.909** | **0.645** | **0.755** | **40** | **4** | **22** |
| **Δ** | **−0.93** | **+0.209** | **+0.419** | **+0.413** | **+26** | **−2** | **−26** |

## Key Takeaways

1. **Gyroscope is critical.** Removing it costs +0.41 event F1 — the model drops from finding 40/62 events to just 14/62. This perfectly validates the feature importance finding (gyro axes = 54% of Gini importance).

2. **Threshold Baseline is unaffected** (identical numbers) — it only uses `accel_mag_max` and `accel_mag_range`, which are both accelerometer features.

3. **Logistic Regression halves in F1** (0.305 → 0.141) without gyro. It loses mostly precision (0.186 → 0.077), meaning without gyro it can't distinguish sit-to-stand from other high-acceleration activities.

4. **For the writeup:** This is your Experiment 3 ablation result. When you validate on Marques (accel-only), expect performance closer to the accel-only row (event F1 ~0.34) rather than the full model (0.76). The gap quantifies exactly how much information the gyroscope carries for this task.

▼ **Explanation**

**The headline for your writeup:** Gyroscope features account for 54% of the Random Forest's feature importance and removing them cuts event-level F1 from 0.755 to 0.341. The model can still detect some sit-to-stands with accelerometer alone (14 of 62 events), but it loses most of its discriminative power. This has a direct practical implication — devices or datasets without gyroscope will see substantially degraded performance.

**This also sets up Experiment 2 nicely.** You now have concrete expectations: Marques (accel-only, elderly) should perform somewhere around the accel-only numbers or worse (since you're also adding domain shift from young to elderly). SisFall (accel+gyro, elderly) should perform somewhere between the full-feature and accel-only numbers (has gyro but has domain shift). If either dataset beats these expectations, that's a positive surprise worth highlighting.

▼ **Phase 6: External Validation — Experiment 2 (Days 10-14)**

## Define Ground Truth Extraction (DO THIS FIRST)

○ **This is the step that takes longer than you think.** Before writing any code, open each external dataset and answer:

- How are sit-to-stand events annotated? (start/end timestamps? protocol structure? video labels?)
- Are labels provided per-sample, per-segment, or per-trial?

- Do you need to derive event times from the protocol description (e.g., "trial D07 is sit-to-stand")?
- **For Marques:** Document exactly how their annotations map to sit-to-stand event start/end times
- **For SisFall:** Document exactly which activity codes correspond to sit-to-stand (e.g., D07, D08) and how to extract event boundaries from the trial structure
- Write this down before coding. If ground truth is ambiguous, note it as a limitation.

## Prepare External Data

- Apply the **exact same** windowing and feature extraction pipeline to Marques and SisFall
- Handle differences:
  - **Sampling rate:** If not 50Hz, resample using `scipy.signal.resample_poly` with anti-aliasing filter. Do NOT use simple interpolation — it introduces aliasing artifacts, especially going from 200Hz → 50Hz.

  - **Sensor availability:** Marques is accel-only → use your accel-only model from Experiment 3 (24 features)
  - **Axis orientation:** May differ from UCI HAPT — this is why your magnitude-based features are valuable (orientation-independent). Document any known differences.
  - **Label format:** Map their labels to your binary scheme using the ground truth extraction rules you defined above

## Train Final Model

- Train RF, logistic regression, and threshold baseline on ALL 30 UCI HAPT subjects
- This is the model you test externally — no retraining on external data

## Run and Report

- Run all three models on Marques data → report precision, recall, F1, PR-AUC
- Run all three models on SisFall data → report precision, recall, F1, PR-AUC
- Apply the same event post-processing (smoothing, min duration, min gap) and report event-level metrics
- Compare against Experiment 1 numbers — how much did performance drop?

- ○ If big drop: discuss domain shift (young vs. elderly, different movement patterns, different sensor characteristics)
- ○ If small drop: the model generalizes well — strong result
- ○ Also report rep-level accuracy on external datasets

▼ **Phase 6.1 - External Validation - Ground Truth**

  ▼ **Results**

Here's the consolidated summary for both datasets:

---

## Dataset 1: Marques (Elderly Accelerometer)

| Property | Value |
|---|---|
| **Location** | `Data/marques/Accelerometer/` + `Data/marques/Results/` |
| **Subjects** | 205 elderly people (IDs 1–205) |
| **Sampling rate** | ~197 Hz (irregular — timestamps in milliseconds) |
| **Sensors** | **Accelerometer only** — 3 axes (X, Y, Z), no gyroscope |
| **Units** | **g's** (magnitude ~9.8 at rest) |
| **Trials per subject** | 6 sit-to-stand repetitions |
| **File format** | Tab-delimited, alternating data/blank lines: `timestamp ax ay az` |
| **Samples per file** | ~19,800 (entire ~100s session in one continuous file) |

**How STS is labeled:** There are **no per-sample labels**. The `Results/results_IDXX.txt` files contain 6 rows of pre-computed metrics per subject (Reaction Time, Total Time, Stand-Up Time, etc.), but these are summary statistics, not segment boundaries. Extracting STS start/end from the continuous accelerometer stream requires either peak-detection or using the timing columns as offsets — this will need careful engineering.

**Implication for our project:** Accel-only means we can only use the accel-only model (24 features). Given the ablation result (event F1 drops from 0.755 → 0.341 without gyro), expect degraded performance. The ~197 Hz rate also differs from UCI HAPT's 50 Hz — we'll need to resample.

---

## Dataset 2: SisFall (Young + Elderly, Accel + Gyro)

| Property | Value |
|---|---|
| **Location** | `Data/sisfall/SA01/` through `SE15/` |
| **Subjects** | 23 young adults (SA01–SA23, ages 19–30) + **15 elderly** (SE01–SE15, ages 60–75) = 38 total |
| **Sampling rate** | **200 Hz** (confirmed: 2400 samples in 12s files) |
| **Sensors** | 2 accelerometers + 1 gyroscope (9 columns total) |
| **Units** | **Raw ADC bits** — must convert using formulas |
| **File naming** | `D07_SA01_R01.txt` = Activity D07, Subject SA01, Trial 01 |
| **File format** | CSV (comma-delimited, semicolon at line end), no header |

**9 Columns:**

| Cols 1–3 | Cols 4–6 | Cols 7–9 |
|---|---|---|
| ADXL345 accel X/Y/Z (±16g, 13-bit) | ITG3200 gyro X/Y/Z (±2000°/s, 16-bit) | MMA8451Q accel X/Y/Z (±8g, 14-bit) |

**Conversion formulas:**

- ADXL345: g = raw × (32/8192) = raw × 0.00391
- ITG3200: °/s = raw × (4000/65536) = raw × 0.0611

- MMA8451Q: `g = raw × (16/16384) = raw × 0.000977`

**Sit-to-stand activities:**

- **D07**: Slowly sit in half-height chair, wait, stand up slowly (190 files, 5 trials/subject)
- **D08**: Quickly sit in half-height chair, wait, stand up quickly (190 files, 5 trials/subject)
- **D09/D10**: Same but low-height chair (185 files each)

**Key detail:** Each trial file is a self-contained 12-second recording of the complete sit→wait→stand sequence. There are **no per-sample labels** within each file — the entire file IS the activity. The sit-to-stand transition happens somewhere in the middle of the 12s window.

**Elderly restrictions:** Subjects SE01–SE15 skipped D06, D13, D18, D19 for safety. Some had additional individual restrictions. SE06 (Judo expert) performed everything.

**Implication for our project:** SisFall has gyro, so we can use the full 48-feature model. The 200 Hz rate needs resampling to 50 Hz. We'll use columns 1–3 (ADXL345 accel) + 4–6 (ITG3200 gyro) to match our 6-channel format. The 15 elderly subjects with D07/D08 trials are the most relevant for cross-dataset validation.

▼

▼ Explanation

**Marques is the harder dataset.** There are no segment labels — just one continuous 100-second file per subject with 6 sit-to-stand reps somewhere inside it. You'd need to build a peak detector to find the rep boundaries in the raw signal, but that's circular — you're trying to validate your rep detector using ground truth that requires a rep detector to create. The pre-computed results files give summary timing but not sample-level boundaries.

**SisFall is much cleaner.** Each trial file is a separate 12-second recording of one sit-to-stand event. The entire file IS the activity. You don't need to find where the sit-to-stand happens — the file boundaries are the labels. D07 (slow) and D08 (fast) with elderly subjects SE01–SE15 are exactly what you need.

Given the timeline pressure (milestone March 6, final March 17), I'd recommend:

**Use SisFall for external validation. Defer Marques or treat it as a limitation.**

Marques would require significant engineering to extract ground truth, and any ground truth you extract would itself be questionable. SisFall gives you clean labels, has both accel and gyro, includes elderly subjects, and has slow vs fast sit-to-stand variants — which directly tests whether your model handles different movement speeds.

▼ **Phase 6.2 - SisFall Prep**

   ▼ **Results**

Everything runs cleanly. Here's the summary:

## SisFall Loader Results

### Dataset Stats

| Property | Value |
|---|---|
| Elderly subjects | 15 (SE01–SE15) |
| D07 (slow STS) trials | 74 (one missing: `D07_SE14_R04.txt`) |
| D08 (fast STS) trials | 75 |
| Total trials | 149 (each = 1 ground-truth event) |
| Samples per trial (50 Hz) | 600 (= 12s × 50 Hz) |
| Windows per trial | 8 (= (600–128)/64 + 1) |
| Total windows | 1,192 |
| Features saved to | `Data/sisfall/features.csv` (1192 × 52) |

### Feature Scale Sanity Check

| Feature | SisFall | UCI HAPT | Match? |
|---|---|---|---|
| `accel_mag_mean` | 1.012 | 1.034 | ✅ Excellent |
| `accel_mag_max` | 1.336 | 1.415 | ✅ Very close |
| `gx_mean` | −0.009 | 0.021 | ✅ Both near zero |
| `gx_std` | 0.373 | 0.347 | ✅ Excellent |
| `gyro_mag_mean` | 0.280 | 0.509 | ⚠️ Lower |

The unit conversions are correct — accel magnitude hovers around 1g at rest, and gyroscope features are in the same order of magnitude as UCI HAPT. The `gyro_mag_mean` is lower in SisFall (0.28 vs 0.51) likely because SisFall trials include long sitting/standing phases where the gyro is near zero, whereas UCI HAPT's windows span a variety of dynamic activities.

## Labeling Strategy

All 1,192 windows are labelled `0`. During external validation, the model will predict on these windows, and event detection checks: **if ≥1 window in a trial is predicted positive, the event is detected.** This gives 149 ground-truth events to evaluate against.

▼ Explanation

The **feature scale sanity check** is making sure that when your model (trained on UCI HAPT data) sees SisFall data, the numbers are in the same ballpark. If they weren't, it would mean something went wrong with unit conversion or resampling.

Here's what we're checking:

`accel_mag_mean` — the average acceleration magnitude in a window. At rest, gravity alone produces ~1g. SisFall shows 1.012, UCI HAPT shows 1.034. Both are right around 1g, confirming the accelerometer units match.

`accel_mag_max` — the peak acceleration in a window. SisFall 1.336 vs UCI HAPT 1.415. Similar range, meaning the movements produce comparable acceleration spikes in both datasets.

`gx_mean` — average gyroscope x-axis rotation. Both near zero (-0.009 vs 0.021), which makes sense — on average, people aren't continuously rotating.

`gx_std` — variability of gyroscope x-axis. SisFall 0.373 vs UCI HAPT 0.347. Very close, meaning rotational movement intensity is comparable.

**gyro_mag_mean** — this one is lower in SisFall (0.28 vs 0.51). Why? Because each SisFall trial is 12 seconds long but the sit-to-stand transition only takes ~2–3 seconds. The rest of the trial is sitting still or standing still, where gyroscope reads near zero. Those quiet periods drag the average down. In UCI HAPT, the windows span all kinds of activities including walking, which keeps gyroscope values higher on average. So the difference is expected — it reflects the data composition, not a conversion error.

**Why this matters:** Your model learned what sit-to-stand "looks like" in UCI HAPT numbers. If SisFall numbers were 10x larger or in completely different units, the model would be useless on them. The fact that they match means the unit conversion and resampling worked correctly, and the model has a fair chance of recognizing sit-to-stand patterns in the new data.

We're comparing. Here's what happened.

**UCI HAPT features** — you computed these in Phase 2. Your pipeline took the raw accelerometer/gyroscope signals, windowed them, and calculated 48 features per window. Those numbers live in `Data/uci_hapt/features.csv`.

**SisFall features** — the script you just ran (`src/sisfall_loader.py`) took SisFall's raw signals, converted the units to match UCI HAPT (g's for accel, rad/s for gyro), resampled from 200Hz to 50Hz, then applied the exact same windowing and feature computation pipeline. Those numbers live in `Data/sisfall/features.csv`.

The sanity check just pulls the mean of a few key columns from each CSV and puts them side by side. If `accel_mag_mean` was 1.01 in one and 9.81 in the other, that would tell you the units don't match (one is in g's, the other in m/s²). If `gx_std` was 0.37 in one and 0.0065 in the other, the gyroscope units are probably wrong (one in rad/s, the other in deg/s).

The fact that the numbers are close means your unit conversion and resampling worked — both datasets are now speaking the same language, so the model trained on UCI HAPT can meaningfully process SisFall data.

▼

▼ **Phase 6.3 - Train, Run, Report**

　　▼ **Results**

Here are the full results from the external validation on SisFall:

---

## External Validation Results: UCI HAPT → SisFall Elderly

## Event-Level Detection (149 trials: 74 D07-slow + 75 D08-fast)

| Model | Overall Recall | D07 (slow) | D08 (fast) | Mean events/ trial |
|---|---|---|---|---|
| **Threshold Baseline (48)** | **0.960** | 0.919 | 1.000 | 0.97 |
| **Logistic Regression (48)** | **0.926** | 0.892 | 0.960 | 1.16 |
| Logistic Regression (24-accel) | 0.148 | 0.095 | 0.200 | 0.18 |
| Random Forest (48) | 0.000 | 0.000 | 0.000 | 0.00 |
| Random Forest (24-accel) | 0.000 | 0.000 | 0.000 | 0.00 |

## Feature-Set Impact

| Model | 48-feat Recall | 24-feat Recall | Δ |
|---|---|---|---|
| Logistic Regression | 0.926 | 0.148 | +0.779 |
| Random Forest | 0.000 | 0.000 | +0.000 |

---

## Key Takeaways

### 1. Does the model generalize to elderly subjects?

- **Yes and no — depends entirely on model complexity.** The Threshold Baseline (96.0% recall) and Logistic Regression (92.6% recall) generalize remarkably well. But the **Random Forest completely fails** — 0/1,192 windows predicted positive, 0/149 trials detected. This is a textbook example of **overfitting to the training domain**. RF learned UCI HAPT-specific decision boundaries so tightly that SisFall's elderly kinematics (different sensor placement, different population, different movement protocol) fall entirely outside them.

### 2. Does D08 (fast STS) get detected more than D07 (slow)?

- **Yes.** Across both working models, D08 consistently outperforms D07:

- TB: D08=100% vs D07=91.9%
- LR: D08=96.0% vs D07=89.2%This makes biomechanical sense — faster sit-to-stand produces larger acceleration peaks and gyroscope angular velocities, making the transition more "salient" to feature thresholds and linear decision boundaries.

**3. How much does accel-only hurt on external data?**

- **Devastating for LR:** dropping from 92.6% → 14.8% recall ($\Delta = -0.779$). Gyroscope features are even more critical for cross-dataset generalization than for within-dataset performance (where the drop was $\Delta = -0.413$). This suggests gyro captures more universal biomechanical signatures of sit-to-stand that transfer across populations.

**4. The bias-variance tradeoff in action:**

- The simplest model (Threshold Baseline) generalizes best. LR is a close second. RF — the best model on UCI HAPT LOSO-CV (F1=0.755) — generalizes worst (F1=0.000). This is the classic result: **models with high internal accuracy don't necessarily transfer**, and simpler models with weaker in-distribution performance can dramatically outperform on out-of-distribution data.

▼

▼ **Explanation 1**

## What each model learned from UCI HAPT

**Threshold Baseline** learned almost nothing specific. It just says: "if there's a big acceleration spike, it's probably sit-to-stand." That's a universal physical fact — standing up from a chair always accelerates your body. It doesn't matter who you are, how old you are, or what phone you're using. Big spike = probably standing up.

**Logistic Regression** learned more. It says: "here's a weighted formula across all 48 features. If the weighted sum is positive, it's sit-to-stand." The formula captures useful patterns like "high acceleration AND high gyroscope rotation together means sit-to-stand." But the formula is a straight line through the feature space — it can't carve out complex, narrow regions. So it captures the general shape of what sit-to-stand looks like, not the exact specifics of UCI HAPT's 30 subjects.

**Random Forest** learned the most. It built 100 decision trees full of very specific rules like: "if `gx_energy` is between 0.08 and 0.15, AND `accel_mag_std` is between 0.07 and 0.12, AND `ax_range` is above 0.63, then sit-to-stand." These narrow ranges were tuned perfectly to how UCI HAPT's young adults move with that specific Samsung phone. The model carved out a very precise, tight region in the 48-dimensional feature space that exactly fits UCI HAPT sit-to-stands.

---

## What happened when they saw SisFall elderly data

Think of it like this. Imagine you're trying to recognize dogs.

**Threshold baseline** learned: "if it has four legs, it's probably a dog." That's crude, but it works on most dogs, regardless of breed, size, or color. It also fires on cats, but it rarely misses a dog.

**Logistic regression** learned: "four legs, fur, medium size, tail up — probably a dog." Better than just "four legs," and it still works across most breeds because the rule is broad enough.

**Random Forest** learned: "four legs, exactly 45cm tall, golden fur, floppy ears, weight between 25–30kg, tail length 28–35cm — dog." This perfectly identifies Golden Retrievers. It never makes a mistake on Golden Retrievers. But show it a Chihuahua and it says "not a dog" because the Chihuahua doesn't fit any of those narrow rules.

In your project:

- UCI HAPT sit-to-stands are the Golden Retrievers — young adults, specific phone, specific protocol
- SisFall elderly sit-to-stands are the Chihuahuas — older people who move differently, different sensors, different setup
- RF learned to recognize UCI HAPT sit-to-stands so precisely that SisFall's sit-to-stands don't match any of its rules
- Not a single one of the 1,192 SisFall windows fell inside RF's narrow decision boundaries

---

## Why this matters for your app

Your app needs to work on **elderly people using various phones**. That's a population and hardware setup very different from UCI HAPT. The results tell you:

A model that's best on your training data (RF, event F1 = 0.755) can completely fail on your target users (RF, event recall = 0.000).

A model that's worse on training data (LR, event F1 = 0.217) can work well on your target users (LR, event recall = 0.926).

This is the **bias-variance tradeoff** in practice. RF has low bias (fits the training data tightly) but high variance (sensitive to differences between training and real-world data). The threshold baseline has high bias (too simple, misses nuance) but low variance (works everywhere). LR sits in between.

For deployment, you might actually want LR or even the threshold baseline — not because they're better models in theory, but because they generalize to the people who will actually use the app.

The Marques dataset on Mendeley is actually the same one you already have — it's the dataset with 40 elderly subjects doing sit-to-stand with accelerometer only and no per-sample labels. The other datasets I found (MotionSense, FoG-STAR) don't have sit-to-stand as a labeled activity.

So realistically, **SisFall is your best external validation dataset**, and you've already run it. There isn't another clean public dataset with labeled sit-to-stand events, elderly subjects, and the right sensor setup.

**How good can you expect Logistic Regression to be?** You already have the answer from your results:

| Setting | LR Event Recall |
|---|---|
| UCI HAPT internal (LOSO-CV) | 39% with strict post-processing |
| SisFall external (elderly) | 92.6% |

The SisFall result is actually surprisingly strong. LR detects 92.6% of sit-to-stand events in elderly subjects it was never trained on. The reason it does better on SisFall than on UCI HAPT internal evaluation is that SisFall trials are isolated 12-second recordings of just sit-to-stand — the model only needs to find one positive window somewhere in the trial. In UCI HAPT, the sit-to-stand is embedded in a continuous stream of many activities, making it harder to distinguish.

**My recommendation: move on to Part 2 (Quality Assessment Layer).** Here's why:

You have a complete and interesting story for Part 1: RF wins internally but fails to generalize, LR generalizes well to elderly subjects, gyroscope is critical, and simpler models transfer better. That's more than enough for the experiments section. Marques can be mentioned as a limitation ("ground truth labels were unavailable for per-sample evaluation"). The time is better spent on Part 2, which is what makes your project unique — the frailty screening indicators that go beyond just counting reps.

## ▼ Explanation 2

The Marques dataset on Mendeley is actually the same one you already have — it's the dataset with 40 elderly subjects doing sit-to-stand with accelerometer only and no per-sample labels. The other datasets I found (MotionSense, FoG-STAR) don't have sit-to-stand as a labeled activity.

So realistically, **SisFall is your best external validation dataset**, and you've already run it. There isn't another clean public dataset with labeled sit-to-stand events, elderly subjects, and the right sensor setup.

**How good can you expect Logistic Regression to be?** You already have the answer from your results:

| Setting | LR Event Recall |
|---|---|
| UCI HAPT internal (LOSO-CV) | 39% with strict post-processing |
| SisFall external (elderly) | 92.6% |

The SisFall result is actually surprisingly strong. LR detects 92.6% of sit-to-stand events in elderly subjects it was never trained on. The reason it does better on SisFall than on UCI HAPT internal evaluation is that SisFall trials are isolated 12-second recordings of just sit-to-stand — the model only needs to find one positive window somewhere in the trial. In UCI HAPT, the sit-to-stand is embedded in a continuous stream of many activities, making it harder to distinguish.

**My recommendation: move on to Part 2 (Quality Assessment Layer).** Here's why:

You have a complete and interesting story for Part 1: RF wins internally but fails to generalize, LR generalizes well to elderly subjects, gyroscope is critical, and simpler models transfer better. That's more than enough for the experiments section. Marques can be mentioned as a limitation ("ground truth labels were unavailable for per-sample evaluation"). The time is better spent on Part 2, which is what makes your project unique — the frailty screening indicators that go beyond just counting reps.

▼

▼ **Phase 7: Writeup & Figures (Days 14-16)**

## Required Tables

○ Experiment 1 results: 3 models × 4 metrics (precision, recall, F1, PR-AUC) — window-level

○ Experiment 1 results: event-level precision, recall, and rep count MAE

○ Experiment 2 results: 3 models × 2 datasets × metrics

○ Experiment 3 results: accel-only vs. full feature set

## Required Figures

○ Confusion matrix for best model (Experiment 1)

○ Feature importance plot (top 10 features from RF)

○ Optional: example raw signal showing sit-to-stand with detected events highlighted

## Key Points to Discuss

○ Which model performed best and why?

○ Did RF justify its complexity over logistic regression?

○ Did the threshold baseline perform surprisingly well or poorly?

○ Did gyroscope features help significantly?

○ How well did models generalize to elderly populations?

○ What were the most common classification errors?

○ Do feature importances align with the literature (Millor, Galán-Mercant)?

○ How did event-level metrics compare to window-level metrics? (This shows whether post-processing rescued noisy window predictions)

▼ **Results**

You already have almost everything on that list. Let me check off what's done:

**Required Tables**

- ✅ Experiment 1 window-level: 3 models × precision, recall, F1, PR-AUC (from LOSO-CV)
- ✅ Experiment 1 event-level: precision, recall, MAE for all 3 models + tuned post-processing (from event detection)
- ✅ Experiment 2: SisFall external validation results (from external validation)
- ✅ Experiment 3: accel-only vs full feature set (from feature ablation)

**Required Figures**

- ✅ Confusion matrix for RF (saved to `Results/confusion_matrix.png`)
- ✅ Feature importance top 10 (saved to `Results/feature_importance.png`)
- ⬜ Optional: raw signal with detected events highlighted (not done yet)

**Key Discussion Points**

- ✅ RF best internally (F1 = 0.554 window, 0.755 event)
- ✅ RF justified over LR internally but failed externally — bias-variance tradeoff
- ✅ Threshold baseline: high recall, useless precision
- ✅ Gyroscope critical: +0.41 event F1
- ✅ Generalization: LR 92.6% recall on elderly, RF 0%
- ✅ False positives are all SITTING windows
- ✅ Feature importance aligns with Millor/Galán-Mercant
- ✅ Event-level metrics much better than window-level for RF with minimal post-processing

All the data exists in your `Results/` CSV files. The tables and discussion points are findings you already have — they'll be written up in the final report, not generated by new code.

The only thing missing is the optional raw signal figure. Everything else just needs to be assembled into the writeup. So yes — you can move to Part 2 now. The writeup (Phase 7) is the last step where you pull everything together into the report format.

---

## Part 2

▼ **Phases**

## Part 2 Checklist: Quality Assessment Layer

**Key clarification:** This layer is entirely rule-based, not ML. It takes the rep boundaries detected by the Part 1 model and extracts clinically meaningful features from the raw sensor data within those boundaries. Each feature is aligned with a specific Fried frailty dimension (weakness, slowness, or exhaustion) and is compared against published reference ranges. Be explicit about this distinction in the writeup — the ML contribution is Part 1; Part 2 demonstrates why accurate rep detection matters.

**Test protocol:** The 30-second chair stand test (30s CST). The user performs as many sit-to-stand repetitions as possible in 30 seconds. This protocol (not the 5-rep STS) is used because: (1) it's what the proposal specifies, (2) it's what Millor et al. used, (3) fatigue/exhaustion indicators (CV, fatigue slope) require 10+ reps to be meaningful — 5 reps is too few for reliable slope estimation.

**Output is educational, not diagnostic.** The app surfaces indicators aligned with published frailty research. It does not diagnose frailty. Interpretation is left to the user or their clinician. Use language like "below expected range for age group" rather than "you are frail." State this explicitly in the writeup.

---

## Phase 1: Extract Per-Rep Features (Days 1-2)

### Prerequisites

- ◯ Part 1 event detection is working — you have start and end timestamps for each detected rep
- ◯ You have access to the raw sensor signals (not just window features) for the corresponding time ranges

### Gravity Removal

- ◯ Before extracting per-rep features, separate gravity from dynamic acceleration
- ◯ Apply a high-pass Butterworth filter (cutoff ~0.3Hz, matching UCI HAPT's preprocessing) to the raw accelerometer signal. This removes the constant gravitational component and leaves only body-movement acceleration.
- ◯ Compute **dynamic acceleration magnitude**: `dynamic_accel_mag = sqrt(filtered_ax² + filtered_ay² + filtered_az²)`
- ◯ This is orientation-independent and represents movement intensity without gravity bias

○ For gyroscope signals, no gravity removal is needed — gyroscope measures rotation directly

## For Each Detected Rep, Go Back to the Raw Signal and Extract:

**Weakness indicators:**

○ **Peak dynamic acceleration magnitude**
  - Within the rep boundary, find `max(dynamic_accel_mag)`
  - This captures how forcefully the person moves during the transition
  - Note: this is NOT "vertical acceleration" — it is total dynamic movement intensity, orientation-independent. Rename from "peak vertical acceleration" throughout to avoid confusion with literature that assumes a defined vertical axis.
  - Reference: Galán-Mercant & Cuesta-Vargas (2013/2014) found frail elderly ~2.7 m/s² vs. ~8.5 m/s² non-frail using iPhone at waist. Their measurements used a defined vertical axis, so direct numeric comparison is approximate — note this as a limitation.

○ **Relative muscle power (adapted Alcázar equation)**
  - Published formula (validated on 5-rep STS): `Power (W/kg) = [0.9 × 9.81 × (height × 0.5 − chair_height)] / (t_5reps / 5 × 0.5)`
  - **Our adaptation for 30s CST:** Use mean time-per-rep from detected reps as the time input: `Power (W/kg) = [0.9 × 9.81 × (height × 0.5 − chair_height)] / (mean_time_per_rep × 0.5)`
  - **State explicitly in the writeup:** "We adapt Alcázar's equation by using mean time-per-rep derived from the 30s CST rather than total 5-rep time. This is an adaptation, not the original validated formula. The adaptation is reasonable because the equation's core relationship — power is inversely proportional to movement time — holds regardless of protocol, but the published normative cut-off points were derived from the 5-rep version."
  - Requires user input: height (m) and chair height (m) — these would be entered in the app
  - Reference: Alcázar et al. (2021) published age- and sex-stratified normative cut-off points (derived from 5-rep STS)
  - Compare output against their thresholds as an approximate reference, noting the protocol difference

**Slowness indicators:**

○ **Time-per-rep**

- Simply: `rep_end_time – rep_start_time` in seconds
- Longer duration = slower movement
- This uses the full rep boundary (sit-to-stand complete cycle), not just the rising phase — reliably segmenting the rising phase alone would require additional signal processing beyond the scope of this project

○ **Peak gyroscope magnitude**
- Within the rep boundary, find `max(gyro_magnitude)` where `gyro_magnitude = sqrt(gx² + gy² + gz²)`
- This captures how quickly the trunk rotates during the transition
- Using magnitude rather than a single gyroscope axis for the same orientation-independence reason as acceleration
- Reference: Millor et al. (2013) found Z-velocity peaks during stand-up were the strongest differentiator across frailty levels (p < 0.001). Their measurements used a defined axis; our magnitude is a proxy.
- Reference: Van Lummel et al. (2013) found lower trunk flexion angular velocity in older adults
- **Fallback if gyroscope unavailable** (e.g., Marques dataset): skip this feature, note it as a limitation, and rely on the remaining five indicators. Do not substitute a fake value.

**Exhaustion / fatigability indicators:**

○ **Coefficient of variation (CV) across all reps in the session**
- For any per-rep metric (e.g., peak dynamic acceleration magnitude), compute: `CV = std(values) / mean(values)`
- High CV = inconsistent performance across the 30 seconds = fatigue or instability
- With 30s CST, you expect 10-20 reps — enough data points for CV to be stable
- Reference: Park et al. (2021) found CV of vertical power range was significantly higher in pre-frail/frail group
- Compute CV for: peak dynamic acceleration magnitude, time-per-rep, peak gyro magnitude, and power

○ **Fatigue slope**
- Plot any per-rep metric (e.g., peak dynamic acceleration magnitude) against rep number (1st rep, 2nd rep, 3rd rep...)
- Fit a simple linear regression: `metric = slope × rep_number + intercept`

- Negative slope = performance declining across the session = fatigability
- Reference: Schwenk/Lindemann found mean velocity decrease of −0.0037 m/s per repetition in older women, with 10/15/20% fatigue thresholds at approximately 8/14/21 reps
- This is the most novel indicator — it captures something invisible to a stopwatch
- Note: the slope value itself is exploratory — no published threshold exists for smartphone-derived fatigue slope during 30s CST. Present it as a trend indicator, not a validated diagnostic.

---

## Phase 2: Compare Against Published Thresholds (Day 3)

### Build a Reference Table

○ Create a table of published thresholds you'll compare against:

| Indicator | Frailty dimension | Threshold / reference value | Source | Match quality |
|---|---|---|---|---|
| Peak dynamic accel magnitude | Weakness | Frail: ~2.7 m/s², Non-frail: ~8.5 m/s² | Galán-Mercant (2013/2014) | Approximate — they used defined vertical axis |
| Relative power (adapted Alcázar) | Weakness | Age/sex-stratified cut-off points | Alcázar et al. (2021) | Approximate — validated on 5-rep STS, we use 30s CST mean |
| Time-per-rep | Slowness | Longer = slower (contextual, compare within session and against session mean) | Millor (2013) | Good — time is protocol-independent |
| Peak gyro magnitude | Slowness | Lower = slower; frail showed significantly lower peaks (p < 0.001) | Millor (2013) | Approximate — they used defined Z-axis |

| CV of per-rep metrics | Exhaustion | Higher = more variable = more fatigued; significantly higher in frail group | Park et al. (2021) | Approximate — they used 5-rep STS with 5 wearable sensors |
|---|---|---|---|---|
| Fatigue slope | Exhaustion | Negative slope = declining; ~−0.0037 m/s per rep | Schwenk/ Lindemann | Exploratory — no validated threshold for this setup |

○ The "Match quality" column is important for the writeup — it shows you understand the limitations of comparing across different protocols and sensor setups

○ Look up the specific Alcázar normative tables — they provide cut-off points by age decade and sex

○ Note: you are NOT building a frailty classifier. You are flagging where individual indicators fall relative to published ranges.

### Implement the Comparison

○ For each indicator, write a simple function that takes the extracted value and returns a status:

- "Within expected range" / "Below expected range" / "Above expected range"
- Base the thresholds on the published values above
- For fatigue slope: "Stable" (slope near zero) / "Declining" (negative slope beyond a threshold)

○ This is just if/else logic — no ML, no training

---

## Phase 3: Build the Three-Tier Output (Day 4)

### Tier 1: Rep Count

○ Total number of detected reps from Part 1

○ This is the standard clinical measure — what a stopwatch gives you

○ Present alongside normative rep count ranges by age/sex (from Jones et al., 1999 — the original 30s CST norms)

### Tier 2: Relative Power Score

○ Adapted Alcázar power value computed from mean time-per-rep, height, and chair height

○ Flag against age/sex normative cut-off points (noting the protocol adaptation)

○ This goes beyond what a stopwatch provides — same rep count can yield different power scores depending on speed

### Tier 3: Movement Quality Flags

○ Present the remaining indicators as individual flags:

- Peak dynamic acceleration magnitude: within range / low (weakness flag)
- Peak gyro magnitude: within range / low (slowness flag) — or "unavailable" if no gyroscope
- CV across reps: within range / high (exhaustion flag)
- Fatigue slope: stable / declining (fatigability flag — exploratory)

○ Each flag is aligned with a specific Fried frailty dimension

○ Each flag has a published reference supporting its clinical relevance (with match quality noted)

### The Early Detection Argument

○ Demonstrate with an example (real or constructed):

- "Subject X completed 12 reps — normal by standard test"
- "But peak dynamic acceleration declined 30% from first to last rep, CV was 0.35 (above Park's threshold for frail group), and adapted power was below Alcázar's age-stratified cut-off"
- "These indicators suggest pre-frailty risk that would be invisible to rep count alone"

○ This is the core narrative of your project — rep count misses pre-frailty, your app catches it

---

## Phase 4: Validate the Pipeline on Your Data (Days 4-5)

### Sanity Checks

- ○ Run the quality assessment pipeline on UCI HAPT subjects using Part 1's detected reps
- ○ Do the extracted values look physiologically plausible?
  - Peak dynamic accel magnitude should be in the range of 1–15 m/s² (not 0.001 or 500)
  - Time-per-rep should be 1–5 seconds (not 0.01 or 30)
  - Power values should be in the range of 0.5–5 W/kg for adults
  - CV should be between 0 and 1 for most subjects (above 1 means extreme variability)
- ○ If values are wildly off, debug the unit conversion, gravity removal, or signal processing

## Compare Across Subjects

- ○ Do subjects who perform fewer reps also tend to have lower peak acceleration and higher CV?
- ○ This isn't a formal experiment — it's a sanity check that the indicators behave as expected
- ○ If a subject with 20 reps has worse quality metrics than a subject with 8 reps, something is probably wrong

## Run on External Data (if available from Part 1)

- ○ Apply the same pipeline to Marques and/or SisFall detected reps
- ○ Marques subjects are elderly — do their quality metrics look different from UCI HAPT (younger)?
  - Expect: lower peak dynamic accel magnitude, lower peak gyro magnitude, higher CV, more negative fatigue slope
- ○ If you see this pattern, it supports the clinical validity of the indicators
- ○ If you don't, discuss why (sensor differences, small sample, protocol differences, etc.)

---

## Phase 5: Writeup (Days 5-6)

### What to Include

- Clear statement that this layer is rule-based, not ML
- Clear statement that the test protocol is 30s CST, not 5-rep STS
- Clear statement that output is educational/screening, not diagnostic
- Table of all six indicators with their frailty dimension alignment, published references, and match quality
- The adapted Alcázar equation written out with variable definitions, with explicit note about protocol adaptation
- Gravity removal method described (high-pass Butterworth, cutoff 0.3Hz)
- Example output showing all three tiers for one or more subjects
- The early detection argument with a concrete example
- Limitations:
  - Thresholds from literature were derived with different protocols (5-rep STS vs. 30s CST), sensor setups, and populations
  - Acceleration magnitude is a proxy for vertical acceleration — direct numeric comparison with literature is approximate
  - Alcázar power formula is adapted, not used in its published validated form
  - No direct validation against Fried phenotype scores (this is future work)
  - Gyroscope-dependent features unavailable with accel-only data
  - Fatigue slope threshold is exploratory — no published cut-off for this setup
  - Small sample sizes in external validation datasets

## Figures

- Per-rep feature plot: e.g., peak dynamic accel magnitude across reps for one subject, showing the fatigue slope line
- Three-tier output summary for one subject (could be a simple table or mockup)
- Optional: comparison of quality metrics between UCI HAPT (younger) and Marques (elderly) subjects

## Connection to Part 1

- Explicitly state: accurate rep segmentation from Part 1 is the prerequisite for all of Part 2
- If Part 1 misses a rep, you lose that rep's quality data
- If Part 1 hallucinates a rep, you get garbage quality values for a non-existent transition

○ This is why rep detection accuracy (especially event-level) matters beyond just counting

---

▼ **Phases 1-2  - Results**

All steps ran successfully. Here's the summary:

---

## Quality Assessment Results

### UCI HAPT — Per-Rep Features (62 reps across 30 subjects)

**Sanity checks — all 62 values within expected ranges:**

| Feature | Min | Max | Mean | Expected |
|---|---|---|---|---|
| Peak dynamic accel (m/s²) | 2.074 | 9.061 | 4.543 | 1–15 ✓ |
| Time per rep (s) | 1.480 | 3.660 | 2.590 | 1–5 ✓ |
| Peak gyro magnitude (rad/s) | 1.006 | 4.922 | 2.070 | 0.5–10 ✓ |
| Relative power (W/kg) | 1.882 | 4.653 | 2.790 | 0.5–5 ✓ |

**Key finding — time vs. acceleration correlation:**

○ Pearson r = **–0.399**, p = 0.029 — subjects with faster reps produce higher peak acceleration. This is the expected biomechanical relationship (faster standing → more forceful push-off).

**Per-subject CV and fatigue slopes** are computed for all 30 subjects but noted as not clinically meaningful (only 2–3 reps each). The pipeline is demonstrated and would be directly applicable to a 30-second chair stand test dataset.

### SisFall Comparison — UCI HAPT (young) vs SisFall (elderly)

The cross-dataset comparison reveals an important caveat:

| Feature | UCI HAPT | SisFall | Notes |
|---|---|---|---|
| Peak dynamic accel | 4.54 m/s² | 8.61 m/s² | ⚠ Unexpected |
| Peak gyro | 2.07 rad/s | 2.75 rad/s | ⚠ Unexpected |
| Time per rep | 2.59 s | 12.00 s | Inflated (full 12s trial) |
| Power | 2.79 W/kg | 0.57 W/kg | Deflated (full 12s trial) |

**Why SisFall elderly show _higher_ peak accel/gyro than expected:** The SisFall D08 (fast sit-to-stand) trials produce very high peaks (12.4 m/s² mean), which pulls the overall SisFall average up. When we look at **D07 (slow) only** — more comparable to UCI HAPT's natural-pace transitions — the picture changes:

| Feature | D07 (slow) | D08 (fast) | Δ |
|---|---|---|---|
| Peak accel | 4.76 m/s² | 12.41 m/s² | +7.65 |
| Peak gyro | 1.69 rad/s | 3.80 rad/s | +2.12 |

**D07 slow STS** (4.76 m/s²) is very close to UCI HAPT (4.54 m/s²), which makes sense — both are natural-pace transitions. D08 fast STS produces ~2.6× higher peaks, confirming the protocol instruction to stand up quickly was followed.

## Files Saved

- `Results/per_rep_features.csv` — 62 rows (one per UCI HAPT rep)
- `Results/subject_quality_summary.csv` — 30 rows (one per subject)
- `Results/sisfall_per_rep_features.csv` — 149 rows (one per SisFall trial)

▼ **Phases 1-2 - Explanation**

Excellent results. The sanity checks all pass, and the findings make sense.

The SisFall comparison is actually more interesting than the expected "elderly = lower" pattern. The D07/D08 split reveals that **protocol matters more than age** in this data — elderly subjects instructed to stand up fast (D08) produce higher peaks than young adults at natural pace (UCI HAPT). Only when you compare similar pace (D07 slow vs UCI HAPT natural) do the values converge. That's a good discussion point: you can't compare raw acceleration values across studies without controlling for movement speed.

The time-per-rep and power values for SisFall are meaningless as expected — the full 12-second trial includes sitting and standing still, so those are artifacts of the trial structure, not real rep durations. Worth a one-sentence note in the writeup.

The negative correlation between time and acceleration ($r = -0.399$, $p = 0.029$) is a nice validation that the pipeline captures real biomechanics — faster standing requires more force, producing higher acceleration peaks.

Go ahead and run Part 2 of 2 — threshold comparison and three-tier output.

**D07, D08, and "natural pace" refer to how fast the person was told to stand up.**

In the SisFall dataset, each activity has a code. D07 means "slowly sit down, wait, then slowly stand up." D08 means "quickly sit down, wait, then quickly stand up." These are instructions given to the participants before each trial — the speed is controlled by the experimenter.

In UCI HAPT, there were no speed instructions. The participants just did their activities naturally — stand up at whatever pace feels normal. That's what I mean by "natural pace."

---

**Now the results make more sense:**

**UCI HAPT (young adults, natural pace):** Average peak acceleration = 4.54 m/s². These are 19–48 year olds standing up at their own comfortable speed. No one told them to go fast or slow.

**SisFall D07 (elderly, told to go slow):** Average peak acceleration = 4.76 m/s². These are 60–75 year olds deliberately standing up slowly. Despite being older, their peak acceleration is almost identical to UCI HAPT's young adults. Why? Because the young adults weren't trying to be fast — their natural pace happens to be similar to an elderly person's slow pace.

**SisFall D08 (elderly, told to go fast):** Average peak acceleration = 12.41 m/s². Same elderly people, but now they're deliberately pushing themselves to stand up as quickly as possible. The acceleration nearly triples compared to their slow trials. Even though they're elderly, the "fast" instruction produces much higher peaks than young adults at natural pace.

---

**What does this mean for your project?**

The original expectation was simple: elderly = lower acceleration, young = higher acceleration. But the data shows it's not that straightforward. How fast someone *chooses* (or is instructed) to move matters more than their age, at least in this comparison. An elderly person standing up fast produces higher peaks than a young person standing up casually.

This is actually a useful finding. In a real 30-second chair stand test, people are told to go as fast as possible. So you'd expect high peaks even from elderly subjects — and the peaks would only start dropping as fatigue sets in over the 30 seconds. That's exactly what the fatigue slope indicator is designed to capture: not the absolute level of acceleration, but whether it declines across reps.

Does that make sense?

▼ **Phase 3-5 - Results**

All three figures look clean and professional. Here's the full summary of what was generated:

---

## `src/three_tier_output.py` — Results Summary

### Classification Distribution (30 UCI HAPT subjects)

| Indicator | Distribution |
|---|---|
| Peak accel | All 30 subjects **Intermediate** (between frail 2.7 and non-frail 8.5 m/s²) |
| Power | 29 Within expected range, **1 Below** (subject 18: 1.92 W/kg) |
| CV accel | 23 Stable, **7 Variable** (exhaustion flag, CV > 0.30) |
| Fatigue slope | 19 Stable, **11 Declining** (fatigability flag) |

This makes sense — UCI HAPT subjects are young healthy adults (19–48 y), so none hit the frail threshold, but the CV and fatigue slope flags still fire for some subjects, demonstrating the pipeline's sensitivity.

### Three-Tier Reports (3 example subjects)

- **Subject 10** (highest accel, 6.99 m/s²): Power ✓, CV ✓, Fatigue ✓ — only the accel "Intermediate" flag since even the strongest subject doesn't reach the 8.5 non-frail cutoff (which was set for vertical-axis-only)
- **Subject 9** (median): All Tier 3 flags stable, clean profile
- **Subject 17** (lowest accel, 2.80 m/s²): **3 flags** — near-frail accel, high CV (0.37), declining fatigue slope (−1.46) — the most concerning profile in the cohort

## Early Detection Example

**Subject 2** — completed 2 reps (looks normal), but quality analysis reveals:

1. Intermediate peak accel (4.04 m/s²)
2. **High CV = 0.470** (exhaustion flag)
3. **Declining fatigue slope = −2.69 m/s²/rep**

→ Pre-frailty risk invisible to rep count alone.

## Figures

1. `fatigue_slope_example.png` — Subject 8 (3 reps) bar chart with regression line (+2.41 slope, r=0.81), frail/non-frail threshold lines
2. `three_tier_example.png` — Clean color-coded clinical report card for Subject 9
3. `young_vs_elderly_comparison.png` — Grouped bars showing UCI HAPT vs SisFall D07 (slow) vs D08 (fast) for peak accel and peak gyro, with threshold lines

   ▾

▾ **Phase 3-5 - Explanation**

Part 2 is complete. Let me flag a few things about these results before you move to the writeup.

**The early detection example with Subject 2 has a problem.** A CV of 0.470 and fatigue slope of -2.69 computed from only 2 reps is not meaningful. With 2 data points, CV is just measuring the difference between those 2 reps, and a regression line through 2 points has no statistical validity. This is the limitation I flagged earlier — UCI HAPT isn't a 30-second CST. For the writeup, I'd recommend either:

1. Being explicit: "We demonstrate the three-tier output using UCI HAPT data. With only 2–3 reps per subject, CV and fatigue slope values are illustrative of the pipeline's functionality, not clinically interpretable. A 30-second CST producing 10–20 reps would yield statistically meaningful exhaustion indicators."

2. Or constructing a hypothetical example with plausible values over 12–15 reps to show what the output *would* look like in a real 30-second test. Label it clearly as constructed.

**Subject 17's profile is your most interesting finding.** Peak acceleration of 2.80 m/s² is right at the frail threshold (2.7 m/s²), despite being a young adult. This could be a genuinely weaker individual, or it could be that they stood up very slowly/gently during the experiment. Either way, it demonstrates that the pipeline catches individual variation that rep count alone would miss.

**The classification distribution makes sense.** All subjects falling in "intermediate" for peak acceleration is expected — Galán-Mercant's thresholds were based on vertical-axis acceleration, while you're using total magnitude, which is always higher. The non-frail cutoff of 8.5 m/s² is hard to reach with magnitude-based measurement. Worth noting this calibration difference in the writeup.

You now have everything needed for the final report. Ready to discuss the writeup structure?

---

## Deliverables

▼ **Completed experiments**
- ○ Experiment 1: LOSO-CV on UCI HAPT (3 models, window + event level)
- ○ Experiment 2: External validation on SisFall elderly
- ○ Experiment 3: Feature ablation (accel-only vs full)
- ○ Feature importance analysis
- ○ Quality assessment pipeline (all 6 indicators, 3-tier output)

▼ **Files**
- ○ All results CSVs in `Results/`
- ○ Confusion matrix, feature importance, fatigue slope, three-tier example, young vs elderly comparison figures
- ○ Sanity check plot

▼ Milestone Report
- ○ **Things you should adjust before submitting:**
- ○ Verify the numbers match your actual outputs (I used the numbers you reported to me, but double-check against your CSVs)

○ You may want to include 1–2 figures (confusion matrix, feature importance plot) if space allows — the 3-page limit excludes references, so you have room

○ Add or remove references as needed — I included the ones most directly relevant

▼ **Final Report**

**What's covered:** All three experiments, the quality assessment layer, the generalization finding (RF overfitting), the feature importance/ablation story, and limitations. The report follows the exact structure the class requires: Introduction, Related Work, Dataset, Methods, Experiments/Results/Discussion, Conclusion/Future Work, References, Contributions.

**Things to verify before submitting:**

○ Double-check all numbers against your actual CSV outputs

○ Reference [4] (Sher et al.) needs the full citation — I had incomplete details from the project outline

○ Reference [3] (Galán-Mercant) — verify this is the correct paper; the 2013 paper on STS acceleration thresholds may be a different publication than the 2014 mRomberg paper

○ If you add figures, reference them in the text (e.g., "see Figure 1")