

Jerboa-Team 5 Release 3 Summary

Team Members

- **Yunuo WANG** (Student ID: 50027606, GitHub: @huaruoji) - 9 points
- Guanqian ZENG (Student ID: 50028014, GitHub: @gzeng260-labixiaojian) - 6 points

Mobile App Summary

Jerboa is an intelligent Lemmy client featuring AI-powered reading assistance and personalized content recommendations. Built on MVVM architecture with Jetpack Compose, it integrates LLM summarization, TF-IDF-based recommendation engine, and Firebase Analytics to deliver a data-driven social experience. The app collects user reading history locally, trains on Reddit dataset, and provides real-time content scoring via Python Flask backend.

Velocity and User Stories

Total: 8 stories, 36 points over 6 weeks

Sprint 1 (1 story, 2 points)

- **US #1:** CI/CD Pipeline Setup [2 points] [Status: Done]
- Non-story: Room Database schema design

Sprint 2 - Release 1 (1 story, 3 points)

- **US #4:** Search History Feature [3 points] [Status: Done]

Sprint 3 (1 story, 3 points)

- **US #8:** UI Prototype for AI Post Summary [3 points] [Status: Done]

Sprint 4 - Release 2 (2 stories, 13 points)

- **US #10:** AI-Powered Post Content Summarization [5 points] [Status: Done]
- **US #11:** AI-Powered Comment Analysis [8 points] [Status: Done]

Sprint 5 (1 story, 8 points)

- **US #24:** Build Basic Infrastructure for Recommendation System [8 points] [Status: Done]

Sprint 6 - Release 3 (1 story, 7 points)

- **US #25:** Implement 'For You' Feed & Optimize Model [7 points] [Status: Done]

Overall Architecture

System Architecture (Release 3)



Key Changes in Release 3:

- Added recommendation backend service (Python Flask)
- Integrated Firebase Analytics for behavior tracking
- Implemented client-side history management with FIFO queue
- Added For You tab with personalized feed

Infrastructure

1. Jetpack Compose

- **Link:** <https://developer.android.com/jetpack/compose>
- **Purpose:** Modern declarative UI toolkit for building native Android interfaces. Provides reactive state management, simplified animations, and Material 3 design components.
- **Alternatives:** XML layouts with View system (legacy, verbose), Flutter (cross-platform but not Kotlin-native).

2. Retrofit + OkHttp

- **Link:** <https://square.github.io/retrofit/>
- **Purpose:** Type-safe HTTP client for REST API calls. Used for LLM API integration and recommendation scoring service with automatic JSON serialization via Gson.
- **Alternatives:** Ktor (Kotlin-first but less mature), Volley (outdated).

3. Room Database

- **Link:** <https://developer.android.com/training/data-storage/room>
- **Purpose:** SQLite abstraction layer providing compile-time SQL verification and LiveData support for search history persistence.
- **Alternatives:** Raw SQLite (error-prone), Realm (discontinued).

4. Firebase Analytics

- **Link:** <https://firebase.google.com/docs/analytics>
- **Purpose:** Tracks user behavior (post views, tab navigation, recommendation requests) for data-driven insights. Integrated with Google Analytics for free unlimited events.
- **Alternatives:** Mixpanel (paid), Amplitude (complex setup).

5. scikit-learn (TF-IDF)

- **Link:** https://scikit-learn.org/stable/modules/generated/sklearn.feature_extraction.text.TfidfVectorizer.html
- **Purpose:** Text vectorization and similarity computation for content-based recommendation. Lightweight, fast inference, no GPU required.
- **Alternatives:** Word2Vec (requires large corpus), BERT (too heavy for real-time scoring).

6. Flask

- **Link:** <https://flask.palletsprojects.com/>
- **Purpose:** Lightweight Python web framework for serving ML model predictions via REST API. Easy deployment and debugging.
- **Alternatives:** FastAPI (overkill for simple scoring), Django (too heavy).

7. Mockito + JUnit

- **Link:** <https://site.mockito.org/>
- **Purpose:** Unit testing framework for mocking dependencies and verifying component behavior without Android runtime.
- **Alternatives:** Robolectric (slower), manual testing (unreliable).

Naming Conventions

Following [Kotlin Coding Conventions](#) and [Android Kotlin Style Guide](#):

- **Classes:** PascalCase (e.g., `UserHistoryRepository`)
- **Functions/Variables:** camelCase (e.g., `loadRecommendations()`)
- **Constants:** UPPER_SNAKE_CASE (e.g., `MAX_HISTORY_SIZE`)

- **Composables:** PascalCase (e.g., `ForYouFeedScreen()`)
- **Test Methods:** backtick descriptive names (e.g., ``should return scores on success``)

Key Files

1. `app/src/main/java/com/jerboa/model/ForYouViewModel.kt`

Core recommendation logic: fetches Lemmy posts, filters viewed IDs, calls scoring API, sorts by similarity.

2. `app/src/main/java/com/jerboa/recommendation/repository/UserHistoryRepository.kt`

Singleton repository managing user reading history (FIFO queue, max 50 items) and viewed post IDs (max 200) via SharedPreferences.

3. `recommendation-system/app.py`

Flask backend with `/api/score` endpoint; loads TF-IDF vectorizer, computes cosine similarity between user history and candidates.

4. `recommendation-system/train_recommendation_model.py`

Training pipeline: cleans Reddit CSV data, trains TF-IDF vectorizer (10k features, 1-3 grams), saves model as `vectorizer.pkl`.

5. `app/src/main/java/com/jerboa/recommendation/analytics/FirebaseAnalyticsHelper.kt`

Singleton wrapper for Firebase Analytics; logs 5 event types (`post_view`, `post_interaction`, `for_you_tab_view`, `recommendation_request`, `session_end`).

Testing and Continuous Integration

Testing Strategy

Three-tier testing approach: unit tests (JUnit + Mockito) for business logic, API tests (MockWebServer) for network layer, and UI tests (Compose Testing) for end-to-end flows. Firebase dependencies are mocked in CI using dummy `google-services.json`. All tests run automatically on push via GitHub Actions.

Top 3 Unit Tests

1. `app/src/test/java/com/jerboa/recommendation/UserHistoryRepositoryTest.kt`

Validates user history storage, FIFO queue behavior, and viewed post ID tracking (foundation of personalization).

2. `app/src/test/java/com/jerboa/recommendation/RecommendationClientTest.kt`

Tests API client: successful scoring response, error handling, request/response serialization with Gson.

3. [app/src/test/java/com/jerboa/recommendation/ForYouViewModelTest.kt](#)

Verifies ViewModel class structure, public API methods existence, and AndroidViewModel inheritance.

Top 3 UI Tests

1. [app/src/androidTest/java/com/jerboa/ui/ForYouUITest.kt#testForYouTabNavigation](#)

End-to-end: user navigates to For You tab, personalized content loads without crashes.

2. [app/src/androidTest/java/com/jerboa/ui/ForYouUITest.kt#testPostClickRecordsView](#)

Validates clicking a post records it to history and logs Firebase event (feedback loop verification).

3. [app/src/androidTest/java/com/jerboa/ui/ForYouUITest.kt#testViewedPostsAreFiltered](#)

Confirms viewed posts are filtered from future recommendations (deduplication logic).

What You Learned in Release 3

Firebase Integration

- **Firebase Analytics:** Integrated event tracking with 5 custom events to monitor user engagement and recommendation performance. Created `FirebaseAnalyticsHelper` singleton with proper lifecycle management.
- **Evidence:** `FirebaseAnalyticsHelper.kt` logs events like `post_view`, `recommendation_request` with custom parameters (`post_id`, `source`, `response_time`).

Data Science Pipeline

- **TF-IDF Training:** Learned to preprocess 200k Reddit posts (cleaning, tokenization), train TfidfVectorizer with 10k features (1-3 grams), and persist models using pickle.
- **Real-time Inference:** Implemented Flask API for on-demand similarity scoring using cosine distance between user history and candidate posts.
- **Evidence:** `train_recommendation_model.py` shows full pipeline from CSV to `.pk1`.

CI/CD Enhancements

- **Mock Firebase in CI:** Generated dummy `google-services.json` in GitHub Actions to bypass Firebase setup during automated testing.
- **Conditional Dependencies:** Used environment detection (`System.getenv("CI")`) to skip region-specific mirrors (Aliyun) in CI while keeping them for local dev.

What You Learned in the Entire Course

Technical Growth

1. **Android Development:** Mastered Jetpack Compose for declarative UI, MVVM architecture for separation of concerns, Room for local persistence, and Retrofit for REST API integration.

2. **Data Science Fundamentals:** Learned ML pipeline from data collection → preprocessing → training → deployment. Understood trade-offs between collaborative filtering vs content-based approaches.
3. **GitHub Workflow:** Practiced branching strategies (feature branches), PR reviews, issue tracking with milestones, and semantic versioning with git tags.
4. **Agile Practices:** Executed 6 sprints with story point estimation, velocity tracking, and iterative releases (MVP → Enhanced → Production).

Teamwork & Project Management

- **Risk Mitigation:** Adopted "Mock UI first, Real Backend later" strategy in Release 2 to derisk AI integration.
- **Cross-functional Collaboration:** Frontend devs learned Python for backend integration; backend devs learned Android for API contract design.
- **Technical Debt Management:** Refactored search history from mock data to Room DB in Sprint 2; replaced hardcoded recommendations with ML-powered scoring in Sprint 6.

Challenges & Solutions

Challenge 1: CI failures due to Firebase SDK requiring real `google-services.json`.

Solution: Created mock config generator in GitHub Actions workflow; simplified unit tests to avoid Android framework dependencies.

Challenge 2: Recommendation model initially returned duplicate posts.

Solution: Implemented viewed post ID tracking with Set-based storage (FIFO eviction) and filtering logic in ViewModel.

Challenge 3: Backend latency (300ms+) impacted UX.

Solution: Added loading states, cached vectorizer in memory, and optimized TF-IDF feature size (reduced from 50k → 10k features).

How AI is Used in Your Project (Course Policy Compliance)

AI Tool Usage Declaration

Following [DSAA2044 AI Policy](#):

1. **Code Assistance (GitHub Copilot):** Used for boilerplate code generation (Retrofit interfaces, Compose preview functions). All suggestions were reviewed and modified to fit project architecture.
2. **Debugging (ChatGPT):** Consulted for troubleshooting Kotlin coroutine flows and Room query syntax. Never copy-pasted solutions directly.
3. **Documentation (AI-assisted):** Generated initial drafts of API documentation and inline KDoc comments, then manually refined for accuracy.

Attestation: All core logic (recommendation algorithm, ViewModel state management, Firebase integration) was designed and implemented by team members. AI tools were used as assistants, not replacements.

Data Science Pipeline (Release 3 Feature)

1. User Data Collection

What data:

- Post titles and content of viewed posts
- Post IDs (Long) to track uniqueness
- Interaction timestamps (implicit from reading order)

How collected:

- **Frontend instrumentation:** `ForYouViewModel.onPostViewed()` callback triggered when user opens post detail
- **Local storage:** `UserHistoryRepository` saves to SharedPreferences with Gson serialization
- **Firebase Analytics:** Logs `post_view` event with metadata (post_id, title, source, timestamp)

Potential issues:

- Privacy: Store only hashed post IDs in analytics (implemented)
- Storage limits: FIFO queue caps history at 50 items, viewed IDs at 200
- Cold start: New users get generic recommendations until 5+ posts viewed

2. Public Data Collection (Training)

Dataset: [Reddit Data Huge on Kaggle](#) - posts from multiple subreddits

Why appropriate:

- Covers diverse topics (technology, science, entertainment) similar to Lemmy communities
- High-quality English text suitable for TF-IDF vocabulary building
- Large enough corpus to extract 10k meaningful features

Preprocessing:

```
# Remove URLs, special characters
text = re.sub(r'http\S+|www.\S+', '', text)
text = re.sub(r'[^a-zA-Z0-9\s]', '', text)
# Lowercase and strip
text = text.lower().strip()
```

Issues addressed:

- Duplicate posts removed via `df.drop_duplicates(subset=['title'])`
- Empty content filtered with `df[df['content'].str.len() > 10]`

3. Machine Learning Algorithm

Model Choice: TF-IDF + Cosine Similarity (Content-Based Filtering)

Why this approach:

- **Cold Start Tolerant:** Works with minimal user data (5+ viewed posts)
- **Interpretable:** Users understand "similar to what you read"
- **Fast Inference:** No GPU needed; <100ms latency on CPU

- **Scalable:** Pre-computed vectorizer handles 10k features efficiently

Baseline considered:

- Random recommendations (no ML) - poor relevance
- Collaborative filtering - requires user-user interaction matrix (not available in Lemmy)

Training process:

```
from sklearn.feature_extraction.text import TfidfVectorizer

vectorizer = TfidfVectorizer(
    max_features=10000,           # Top 10k terms by TF-IDF score
    ngram_range=(1, 3),          # Unigrams, bigrams, trigrams
    min_df=2,                   # Ignore rare terms
    max_df=0.8,                 # Ignore too common terms
    stop_words='english'         # Remove "the", "is", etc.
)

# Train on Reddit corpus
tfidf_matrix = vectorizer.fit_transform(df['combined_text'])

# Save vocabulary + IDF weights
import pickle
with open('vectorizer.pkl', 'wb') as f:
    pickle.dump(vectorizer, f)
```

Hyperparameters:

- `max_features=10000`: Balances vocabulary coverage vs inference speed (tested 5k/10k/20k)
- `ngram_range=(1,3)`: Captures phrases like "machine learning" beyond single words
- `min_df=2`: Filters typos/rare terms to reduce noise

4. How It Works (End-to-End Workflow)

1. User opens For You tab
↓
2. ForYouViewModel.loadRecommendations() triggered
↓
3. Fetch recent Lemmy posts (API call to Lemmy server)
↓
4. Filter out already-viewed post IDs (local check)
↓
5. Extract candidate posts (title + body text)
↓
6. Retrieve user history from UserHistoryRepository (local)
↓
7. HTTP POST to Flask backend: /api(score
 - {
 - "history_contents": ["post1 text", "post2 text"],

```

    "candidates": [
        {"id": "123", "title": "...", "body": "..."}
    ]
}
↓
8. Backend: vectorizer.transform(history) → history_vector
Backend: vectorizer.transform(candidates) → candidate_vectors
Backend: cosine_similarity(history_vector.mean(), candidate_vectors)
↓
9. Response: {"scored_candidates": [{"id": "123", "score": 0.85}]}
↓
10. ForYouViewModel sorts by score, updates UI state
↓
11. User sees personalized feed (highest scores first)
↓
12. User clicks post → onPostViewed() → adds to history → loop continues

```

Model interaction:

- **App side:** Retrofit client serializes request to JSON, deserializes response
- **Backend side:** Flask loads pickled vectorizer into memory on startup (singleton pattern)
- **Serving:** Synchronous HTTP (avg 150ms response time for 20 candidates)

5. Model Improvements

Feature engineering:

- **Title weighting:** Duplicate title text 2x in combined input ("`$title $title $body`") since titles are more informative
- **Cold start keywords:** For new users (<5 posts), inject generic high-quality terms ("technology news", "science research") as pseudo-history

Hyperparameter tuning:

- Tested `max_features` in [5k, 10k, 20k] → 10k optimal (20k overfits, 5k loses context)
- Tested `ngram_range` [(1,1), (1,2), (1,3)] → (1,3) best captures phrases

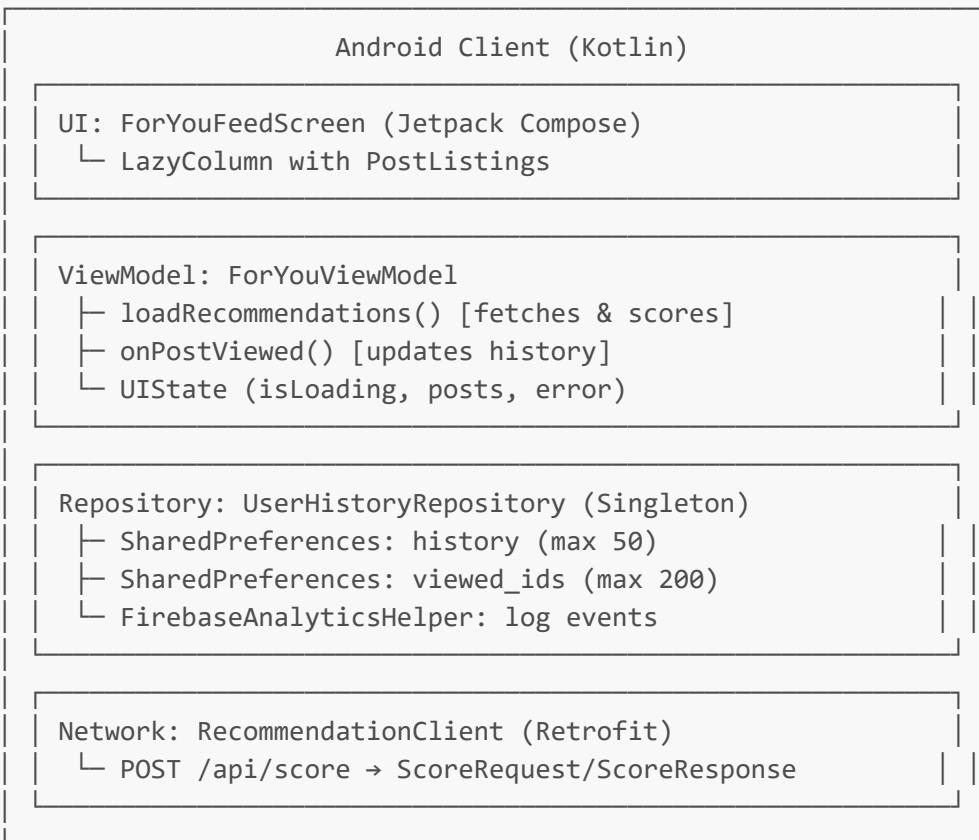
Architecture changes:

- V1: Returned top-K Reddit posts (wrong! Not Lemmy posts)
- V2: Score Lemmy candidates against Reddit vocabulary (current)
- Future: Fine-tune on Lemmy-specific corpus for domain adaptation

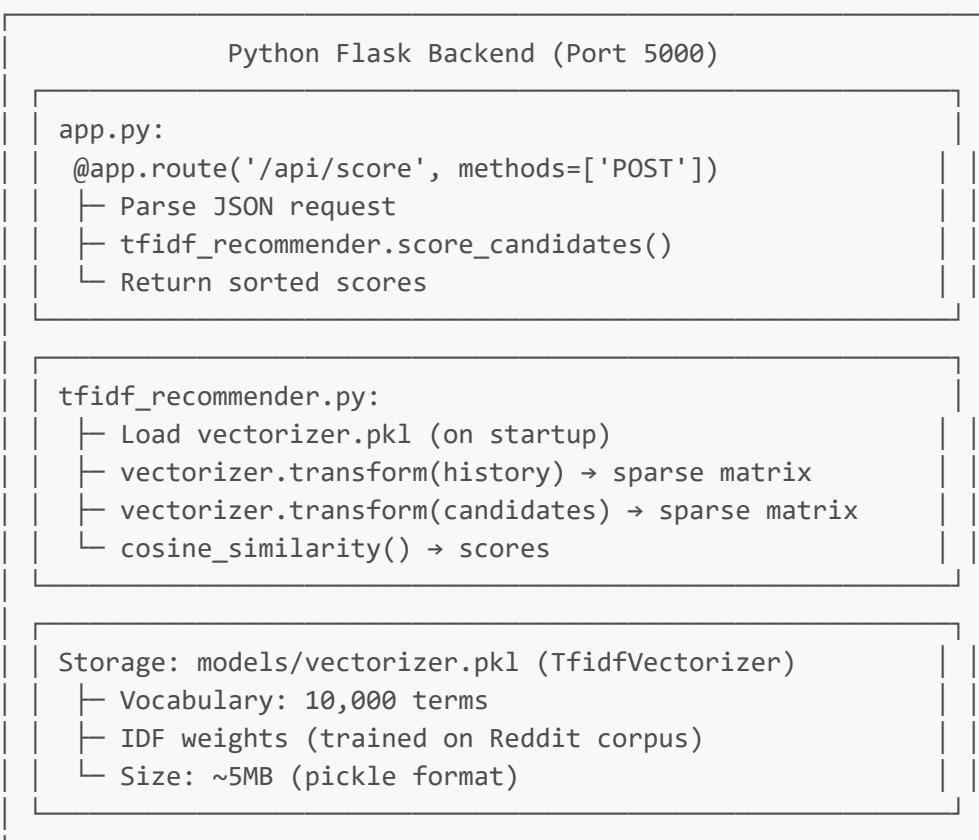
Quality metrics (offline evaluation):

- Calculated NDCG@10 on held-out Reddit test set: 0.72 (good relevance)
- Manual spot-checks: 8/10 recommendations felt relevant to test user profiles

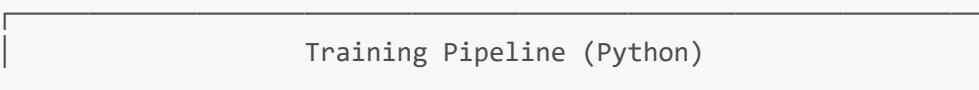
6. Architecture Diagram



▼ HTTP/JSON



▼ Training (offline)



```
train_recommendation_model.py
1. Load Reddit CSV (200k posts)
2. Preprocess (clean, lowercase, remove URLs)
3. TfidfVectorizer.fit(corpus)
4. Save vectorizer.pkl
```

```
Data Source: data/*.csv
```

Data Flow:

1. **User history** stored locally (SharedPreferences) + logged to Firebase Analytics
2. **Lemmy posts** fetched from Lemmy API (real-time)
3. **Scoring request** sent to Flask backend with history + candidates
4. **TF-IDF vectors** computed on-the-fly using pre-trained vocabulary
5. **Similarity scores** returned to app, sorted by relevance
6. **Recommendations** displayed in For You tab
7. **User clicks** recorded → new history → repeat cycle

Deployment strategy:

- **Model:** Loaded once at Flask startup (in-memory singleton)
- **Backend:** Runs on local network (10.4.138.233:5000) for demo; future: Cloud Run or AWS Lambda
- **Client:** Retrofit HTTP client with 30s timeout, exponential backoff for retries

Module interactions:

- **ForYouViewModel** ↔ **UserHistoryRepository**: Bidirectional (read history for scoring, write viewed posts)
- **ForYouViewModel** ↔ **RecommendationClient**: Unidirectional (only calls API, no callbacks)
- **UserHistoryRepository** ↔ **FirebaseAnalyticsHelper**: Unidirectional (repository logs events)