

Multi-Objective Recommender Systems for Enterprise Social Networks:

Adaptive Target Optimization Under Organizational Heterogeneity

Technical Report

Anjan Goswami

Abstract

Enterprise social networks present a recommendation problem fundamentally different from consumer platforms: engagement semantics vary across organizations, interaction types carry different business value depending on context, and cold-start is the norm rather than the exception. This technical report describes the design and deployment of a multi-entity recommender system for Salesforce Community Cloud, serving recommendations across four entity types—individuals, groups, files, and articles—to millions of users across thousands of organizations. The system replaced a heuristic rule-based engine with a unified matrix factorization framework that projects heterogeneous entities into a shared latent space, scored by a multi-objective function combining behavioral signals (likes, shares, comments, follows, reads) with learned weights. The central technical contribution is an adaptive target optimization framework: default objective weights are learned from aggregate platform data, then adjusted per-organization through a Bayesian feedback mechanism calibrated by human evaluation. The system improved engagement by over 20% while providing organizational administrators with transparency, override capabilities, and cost-efficient nightly batch execution.

1 The Problem: Recommendations in Enterprise Social Networks

Salesforce Community Cloud provides organizations with branded online communities—spaces where employees, customers, and partners interact through posts, files, groups, and knowledge articles. The platform’s value depends on engagement: users must discover relevant people to follow, groups to join, files to read, and articles to consult. Without effective recommendations, communities suffer from low engagement, information silos, and the “empty room” problem where new users see no activity relevant to them.

The existing recommendation system was **rule-based**: simple keyword matching on user profiles and content metadata, augmented by recency and popularity heuristics. It suffered from three fundamental limitations:

- **No semantic understanding:** Recommendations were driven by surface-level text overlap between user profiles and content. Two users working on the same problem but using different vocabulary were invisible to each other. A file titled “Q3 Revenue Analysis” would not be recommended to a user whose profile mentioned “sales forecasting.”
- **No behavioral learning:** The system ignored the richest signal available—how users actually interacted with content. Likes, shares, comments, and reading patterns were not incorporated into recommendations.
- **No cross-entity reasoning:** Person recommendations, group recommendations, file recommendations, and article recommendations operated as four independent systems with no shared

intelligence. A user’s group subscriptions revealed nothing about which files they might want; their reading patterns informed nothing about which people they should follow.

The result was low recommendation engagement: click-through rates, follow rates, and content interaction rates were all well below platform benchmarks for social recommendation systems.

1.1 The Heterogeneity Challenge

The deeper problem—and the one that makes enterprise recommendation fundamentally harder than consumer recommendation—is **organizational heterogeneity**. A customer support community values different engagement patterns than an internal R&D collaboration space. In a support community, article reads and file downloads are the primary success metric; in a collaboration space, comments and shares indicate value. A recommendation system that optimizes for a single engagement objective will be wrong for most organizations.

This rules out the standard approach of training a single model on aggregate data and deploying it globally. The target function itself must adapt to each organization’s engagement semantics.

2 System Architecture

2.1 Entity Types and Interaction Signals

The system recommends four entity types to individual users, each with distinct interaction signals:

Recommendation Type	Primary Signals	Intent Captured
Individual → Individual	Like, share, comment, follow	Professional affinity
Group → Individual	Subscription, post activity	Topical interest
File → Individual	Like, share, comment, read	Information need
Article → Individual	Like, share, comment, read	Knowledge seeking

Table 1: Entity types and associated behavioral signals.

The signal space is heterogeneous: *follow* and *subscribe* are explicit affinity signals with low frequency but high precision; *like* is a lightweight endorsement; *share* implies the content is valuable enough to propagate; *comment* implies engagement depth; *read* (for files and articles) captures consumption without explicit endorsement. Each signal carries different information about user preference, and their relative importance varies by entity type and organizational context.

2.2 Unified Latent Space via Matrix Factorization

Rather than building four independent recommendation engines, we constructed a **unified latent space** that embeds all entity types—users, groups, files, and articles—into a shared k -dimensional representation. This allows cross-entity reasoning: a user’s interactions with files inform their group recommendations, and their group memberships inform people recommendations.

Let \mathcal{U} denote the set of users, \mathcal{G} the set of groups, \mathcal{F} the set of files, and \mathcal{A} the set of articles within an organization. We learn latent factor matrices:

$$\mathbf{U} \in \mathbb{R}^{|\mathcal{U}| \times k} \quad (\text{user embeddings}) \quad (1)$$

$$\mathbf{G} \in \mathbb{R}^{|\mathcal{G}| \times k} \quad (\text{group embeddings}) \quad (2)$$

$$\mathbf{F} \in \mathbb{R}^{|\mathcal{F}| \times k} \quad (\text{file embeddings}) \quad (3)$$

$$\mathbf{A} \in \mathbb{R}^{|\mathcal{A}| \times k} \quad (\text{article embeddings}) \quad (4)$$

For a user u and candidate entity $e \in \mathcal{G} \cup \mathcal{F} \cup \mathcal{A} \cup \mathcal{U}$, the base affinity score is the inner product in the latent space:

$$s_{\text{base}}(u, e) = \mathbf{u}_u^\top \mathbf{v}_e$$

where \mathbf{u}_u is the u -th row of \mathbf{U} and \mathbf{v}_e is the corresponding row of \mathbf{G} , \mathbf{F} , \mathbf{A} , or \mathbf{U} depending on entity type.

2.3 Topical Vocabulary Processing

Enterprise communities use specialized vocabulary that varies dramatically across organizations—a pharmaceutical company’s community discusses “Phase III trials” and “FDA submissions” while a technology company discusses “sprint velocity” and “CI/CD pipelines.” Standard NLP preprocessing (generic stopwords, universal tokenization) discards domain-critical terms or treats them as noise.

We built an **organization-specific vocabulary extraction pipeline** that:

- Extracted a topical vocabulary per organization from the corpus of posts, files, and articles using TF-IDF with organization-level IDF computed against cross-organization background frequencies. Terms with high within-org frequency but low cross-org frequency were flagged as domain-specific.
- Constructed topic representations for each entity by aggregating over the organization-specific vocabulary, producing dense topical features that supplemented the collaborative filtering signals from the matrix factorization.
- Updated vocabulary quarterly to capture evolving organizational language (new product names, project codes, emerging initiatives).

These topical features were incorporated as side information in the factorization, regularizing the latent factors toward topically coherent recommendations—particularly important for cold-start entities with sparse behavioral data.

3 Multi-Objective Target Design

3.1 The Composite Engagement Score

The central design problem is defining what “good recommendation” means when multiple interaction types exist and their relative value varies by context. We formalize this as a **weighted composite engagement score**.

For a recommended entity e shown to user u , let ℓ_{ue} , h_{ue} , c_{ue} , f_{ue} , and r_{ue} denote binary indicators for whether user u performed a like, share, comment, follow/subscribe, or read action on entity e , respectively. The composite engagement target is:

$$y_{ue} = \alpha \cdot \ell_{ue} + \beta \cdot h_{ue} + \gamma \cdot c_{ue} + \delta \cdot f_{ue} + \rho \cdot r_{ue}$$

where $\theta = (\alpha, \beta, \gamma, \delta, \rho)$ is the weight vector over engagement types. Not all signals apply to all entity types: f_{ue} (follow/subscribe) applies to individuals and groups; r_{ue} (read) applies to files and articles. The weight vector is entity-type-specific:

$$\theta_{\text{person}} = (\alpha, \beta, \gamma, \delta_{\text{follow}}, 0) \tag{5}$$

$$\theta_{\text{group}} = (\alpha, \beta, \gamma, \delta_{\text{subscribe}}, 0) \tag{6}$$

$$\theta_{\text{file}} = (\alpha, \beta, \gamma, 0, \rho_{\text{read}}) \tag{7}$$

$$\theta_{\text{article}} = (\alpha, \beta, \gamma, 0, \rho_{\text{read}}) \tag{8}$$

3.2 The Scoring Function

The full scoring function for user u and entity e of type t combines the latent affinity with the multi-objective target:

$$\hat{y}_{ue} = \sigma(\mathbf{u}_u^\top \mathbf{v}_e + \mathbf{x}_{ue}^\top \mathbf{w}_t)$$

where \mathbf{x}_{ue} is a feature vector incorporating topical similarity, recency, popularity, and organizational context features, \mathbf{w}_t is an entity-type-specific weight vector, and σ is the sigmoid function. The model is trained to minimize the weighted loss:

$$\mathcal{L} = \sum_{(u,e) \in \mathcal{D}} \sum_{s \in \mathcal{S}_t} \theta_s^{(t)} \cdot L\left(\hat{y}_{ue}^{(s)}, y_{ue}^{(s)}\right) + \lambda (\|\mathbf{U}\|_F^2 + \|\mathbf{V}\|_F^2)$$

where \mathcal{S}_t is the set of applicable signals for entity type t , L is the binary cross-entropy loss for each signal, and λ controls regularization. The key insight is that the θ weights control not just how engagement is measured after the fact, but *what the model learns to optimize*—they shape the loss landscape itself.

3.3 Finding Optimal Target Weights

The weight vector θ is not a hyperparameter to be tuned on a validation set in the traditional sense. It encodes a *value judgment* about which engagement types matter, and this judgment varies by organization. The optimization problem is:

$$\theta^* = \arg \max_{\theta} \mathbb{E}[\text{Engagement}(\theta) | \text{org}] \quad \text{subject to} \quad \theta \geq 0, \quad \|\theta\|_1 = 1$$

We decompose this into two stages:

Stage 1: Platform-default weights from Salesforce data. We had direct access to Salesforce's own Community Cloud deployment—the internal instance used by Salesforce employees. Using this data, we trained the full system and conducted a grid search over θ values, evaluating each configuration against a human-labeled relevance set (described in Section 5). The optimal weights on Salesforce data became the **platform defaults**:

$$\theta_{\text{default}} = \arg \max_{\theta} \text{NDCG}@k(\text{Rank}(\theta), \text{HumanLabels})$$

Stage 2: Per-organization adaptation from behavioral feedback. For customer organizations, we could not obtain human relevance judgments at scale. Instead, we used observed behavioral feedback to adapt θ away from the default. The adaptation follows a Bayesian updating rule:

$$\theta_{\text{org}}^{(t+1)} = (1 - \eta) \cdot \theta_{\text{org}}^{(t)} + \eta \cdot \hat{\theta}_{\text{observed}}^{(t)}$$

where $\hat{\theta}_{\text{observed}}^{(t)}$ is estimated from the organization's recent engagement distribution:

$$\hat{\theta}_s \propto \frac{n_s}{\sum_{s'} n_{s'}} \cdot q_s$$

Here n_s is the count of signal type s observed in the most recent feedback window, and q_s is a quality multiplier derived from the signal's correlation with sustained engagement (e.g., a comment that leads to a follow-on interaction receives higher q_s than a solitary like). The learning rate η controls how aggressively the system departs from platform defaults—set conservatively (0.1–0.2) to prevent overfitting to short-term behavioral noise.

This two-stage approach resolves a fundamental data availability problem. Most customer organizations have insufficient behavioral data to estimate θ from scratch—they face a cold-start problem not just for items but for the *objective function itself*. The platform defaults provide an informed prior; behavioral feedback provides the likelihood; the adapted weights are the posterior.

4 Operational Design

4.1 Batch Execution and Cost Management

The recommendation pipeline executed as a **nightly batch process**. This was a deliberate architectural decision balancing recommendation freshness against computational cost. Real-time recommendation generation was infeasible at the scale of thousands of organizations, each with their own embedding space and parameter configuration.

The nightly cadence provided a natural interval for:

- **Embedding recomputation:** Updating the matrix factorization with the previous day’s interaction data. Incremental updates rather than full retraining kept computation within the batch window.
- **Parameter adaptation:** Updating θ_{org} based on accumulated behavioral feedback since the last cycle.
- **Recommendation materialization:** Pre-computing top- N recommendations for each user across all four entity types, stored for low-latency serving during the day.

The batch interval was a tunable parameter with cost implications. Shorter intervals (e.g., every 6 hours) improved freshness but increased compute costs proportionally. We found that nightly execution provided sufficient freshness for enterprise communities, where content velocity is orders of magnitude lower than consumer social networks.

4.2 Administrative Controls and Override Mechanisms

Enterprise deployments require organizational administrators to maintain control over recommendations. We built a three-layer control system:

- **Monitoring dashboard:** Real-time visibility into recommendation performance—engagement rates by entity type, distribution of recommended entities, and trend analysis. Administrators could see which recommendations were driving engagement and which were being ignored.
- **Promotional slots:** Administrators could pin specific entities (a new group, an important article, a featured employee) to the top of recommendation lists for a configurable duration. Promotional items occupied dedicated slots that did not displace algorithmic recommendations, preventing the promotion mechanism from degrading recommendation quality.
- **Override and exclusion rules:** Administrators could exclude specific entities from recommendations (e.g., deprecated files, inactive groups) or adjust the relative weight of entity types in the recommendation mix. These overrides took effect in the next batch cycle.

This design reflected a core principle: **in enterprise settings, the algorithm serves the administrator, not the other way around**. Organizational context that the model cannot capture—an upcoming reorganization, a compliance-sensitive topic, a strategic initiative—must be expressible through administrative controls.

5 Validation: Human Evaluation and Calibration

5.1 Labeling Protocol

We recruited human evaluators through a vendor in India to assess recommendation quality. The evaluation served two purposes: validating overall system quality and calibrating the target weight optimization.

Each evaluator rated recommendation lists on a graded relevance scale (0–4) across two dimensions:

- **Topical relevance:** Is this entity related to the user’s apparent interests and activity?
- **Actionability:** Would a reasonable user be likely to interact with this recommendation?

Evaluators were provided with the user’s profile, recent activity history, and organizational context. Each recommendation list was rated by three independent evaluators, with inter-annotator agreement measured by Krippendorff’s α .

5.2 Calibrating Target Weights with Human Judgments

The human evaluation data served as the ground truth for Stage 1 of the target weight optimization. For each candidate weight vector θ , we generated recommendation lists, measured their NDCG against human relevance ratings, and selected the θ that maximized human-judged quality.

This calibration revealed important structure. Across entity types:

- **Comments carried the highest quality signal:** Recommendations optimized with high γ (comment weight) consistently scored highest on human-judged relevance and actionability. Comments are costly signals—they require cognitive effort—and their presence strongly indicated genuine engagement.
- **Likes were necessary but not sufficient:** Optimizing primarily for likes (α dominant) produced popular but generic recommendations. Likes provided coverage (dense signal, many observations) but low precision for individual relevance.
- **Reads were essential for files and articles:** For content entity types, the read signal (ρ) was the strongest single predictor of human-judged quality. Users who read a file found it useful; users who liked a file may have endorsed it socially without reading it.
- **Shares were the strongest cross-entity signal:** A user who shared an article was likely to be interested in the article’s author (person recommendation), the group where it was posted (group recommendation), and related files. Share events, though sparse, produced the highest cross-entity information transfer.

These findings informed the default weights, which were not equal: $\gamma > \beta > \rho > \delta > \alpha$ for the general case, with entity-type-specific adjustments (e.g., ρ elevated for files and articles).

6 Results

Metric	Result
Overall engagement lift	>20% across recommendation surfaces
Recommendation click-through rate	Significant improvement over heuristic baseline
Cross-entity discovery	Users exposed to multi-entity recommendations engaged with entity types they had not previously interacted with
Cold-start performance	New users received relevant recommendations within first session via topical features
Administrative adoption	Dashboard actively used by community managers

Table 2: System performance relative to heuristic baseline.

The engagement lift was not uniformly distributed. Organizations with higher baseline activity saw larger absolute improvements, because the behavioral feedback loop had more signal to learn from. Organizations with sparse activity relied more heavily on the platform defaults and topical features, with improvements concentrated in the content recommendation types (files and articles) where topical vocabulary matching provided strong cold-start coverage.

7 Generalizable Principles

1. In multi-stakeholder platforms, the target function is the product. The most consequential design decision was not the model architecture—it was the composite engagement score and the mechanism for adapting it per organization. In consumer recommendation, the platform defines success (e.g., watch time, purchases). In enterprise recommendation, each customer organization defines success differently. The system must optimize for a target that it does not fully know at deployment time.

2. Matrix factorization across heterogeneous entity types creates cross-entity intelligence. The unified latent space was the key architectural decision. Four independent recommendation engines would have produced reasonable per-entity recommendations. The shared space produced *cross-entity discovery*—recommendations that connected users to entity types they had not previously engaged with, based on latent affinity inferred from their interactions with other entity types.

3. Default parameters from a reference deployment solve the cold-start-of-objectives problem. When you cannot collect ground truth from every customer, use your own deployment as the reference. Salesforce's internal community provided the data to calibrate default target weights. This is a general pattern for enterprise ML: the vendor's own usage of the product generates the prior that bootstraps customer deployments.

4. Human evaluation calibrates what behavioral signals cannot. Behavioral signals tell you what users *did*; human evaluation tells you what they *should have seen*. The labeling study revealed that comment-weighted targets produced higher-quality recommendations than like-weighted targets—a finding that could not have been derived from behavioral data alone, because the behavioral data reflected recommendations generated under a different objective.

5. Administrative controls are not a concession—they are a feature. Enterprise customers will not adopt a recommendation system they cannot understand and override. The monitoring dashboard and promotional slots were not add-ons; they were requirements for organizational trust. Algorithmic recommendations that administrators cannot explain or adjust are a liability in enterprise settings.

6. Batch execution is a valid architecture when content velocity is low. The instinct in recommendation systems is to pursue real-time scoring. For enterprise social networks with content velocities orders of magnitude below consumer platforms, nightly batch execution provides sufficient freshness at dramatically lower infrastructure cost. The savings in compute were reinvested in per-organization model adaptation—a higher-value use of the same budget.

7. Behavioral adaptation of objective weights closes the loop that human evaluation opens. Human evaluation sets the prior; behavioral feedback provides the likelihood. The exponential moving average update rule for per-organization weights ensures that the system converges toward each organization's revealed preferences over time, without requiring ongoing human evaluation for every customer.