

the music tracks (total of 1.5M artists). We then processed the MusicBrainz DB⁵ to obtain artist information for each of each track to the corresponding artist. To this end we derived To obtain a more granular feature space, we decided to map each track is labeled with a score for a set of 100 music genres. 221M listening events and around 4.6M music tracks, where covering a period of 2-years (May 2003-May 2014), containing the full “scorpled” listening history of a set of 44.124 users, track logging process is called scorpling. Our dataset contains using a collaborative filtering algorithm. This automated job charts and musical recommendations are calculated. All songs played are added to a log from which personalized for users identified by the platform), or the user’s “friends”, based on the user’s profile, its “musical neighbors” (i.e., sim- consist of uninteracted audio streams of individual tracks cial client application, or with the web player. Radio stations tending to the Last.fm Radio service, either with Last.fm off- player supporting the Last.fm Audiocoroller plugin, or by lis- listening to their personal music collection with a music system Last.fm. Users update their profiles in multiple ways: data from the popular online (social) music recommendation data we considered a dataset from [40] containing temporal

Last.fm Dataset. For an evaluation of our algorithm on real with probability $\frac{10}{I}$, removing any non-zero entry immediately after insertion 1-entries of each row. We introduced deletions by randomly the insertion-only case, the stream consists of the sequence of pily $1/3$ and added a new item with probability $\frac{10}{I} \cdot \frac{3}{I} = \frac{30}{I}$. In (0.42, 0.22) we deleted an item contained in row i with proba-

No. of pairs similarity	10 0.4	1 0.42	3 0.2	0 ≥ 0.22
No. of pairs similarity	1e454e 0.5	2822 0.52	5ee 0.3	13 0.32
No. of pairs similarity	ee435110 0.0	32e42482 0.02	150312e2 0.1	1e38112 0.12

Pairs in Last.fm Dataset
Distribution of Exact Similarity Values for
TABLE 1

the range specified by line 2 of Algorithm 5. This allows to ing set cardinality and similarity threshold are such that k is tures for $\mathcal{L}_{(i)}^k$ only if the bucket is sensitive, i.e., its correspond- in case of deletions of an item λ , we recompute a set of signa- is the selective recomputation of signatures in case of deletions. further optimization that we implemented in D22 Proactive, tures only for the two compressed bucket sets \mathcal{L}_i^k and \mathcal{L}_i^{0*} . A the full user profile, D22 Proactive has to recompute signa- recompute signatures, yet while Vanilla-MH has to do so for mmm. Let $k = \text{isp}(\nu(\lambda))$. In case of deletions, both will have to element and updated in case such value is the new mini- ment is added, all hash-functions are evaluated on the new ture and Vanilla-MH behave in the same way. When an ele- online. When inserting item λ added to set i , both D22 Proac-

Let us now focus on the algorithms that update signatures ing and storing the signatures of sets. responsiveness and additional space required for comput- depends on the use case, with a trade-off between delay

The choice of the first or the second implementation most recent change from the stream).

every update (that is, after line 2 of Algorithm 1, reflecting the Proactive), instead maintains a set of fingerprints online, with The second, called DynamicSketch Proactive (or simply D22 of Algorithm 1 and computes fingerprints only at query time. previously called DynamicSketch (D22) maintains the sketches

We tested two versions of our algorithm. The first version functions (that is also our signature scheme). sketches are computed online using 5-wise independent hash “vanilla” LSH scheme (later Vanilla-MH), where profile compare our approach with an online implementation of a from our synthetic dataset. As a comparable benchmark, we tion-only stream, and a fully dynamic stream, both obtained to various dataset sizes, in two different scenarios, an inser- synthetic dataset, to understand its performance with respect We evaluated the running time of our algorithm using the

6.1 Performance Evaluation

larity values for all pairs the Last.fm dataset. of 0.5M entries. Table 1 shows the distribution of exact simi- $n = 12K$ users (sets), $|U| = 380K$ (items) and a stream length