

Selection of audio features for music emotion recognition using production music

Chris Baume¹, György Fazekas², Mathieu Barthet², David Marston¹ and Mark Sandler²

(1) **BBC** R&D

(2) **centre for digital music**

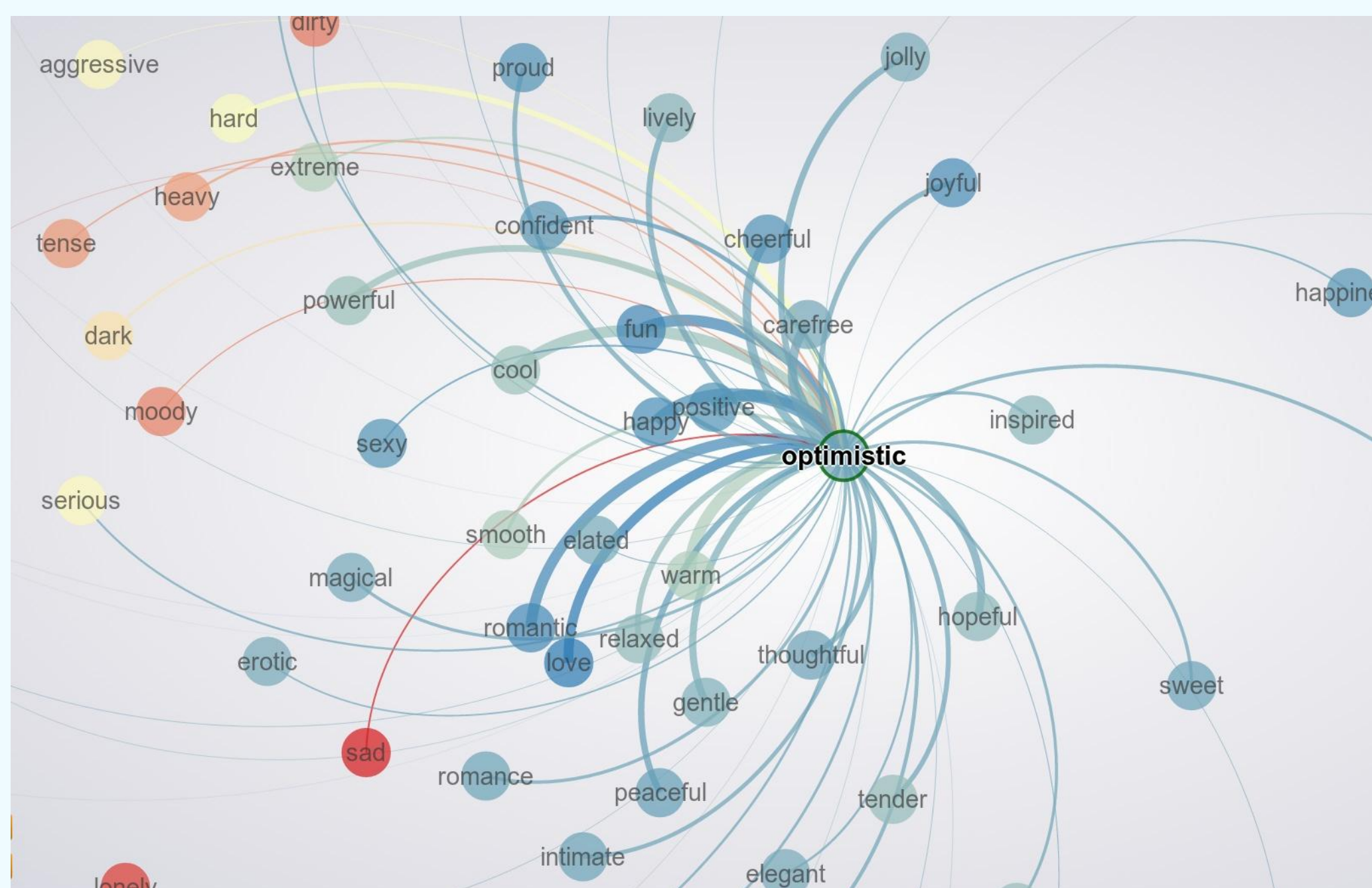
Which combination of audio features gives the best performance for music emotion recognition?

Our experiment used **support vector regressors** to map between various combinations of audio features and a five dimensional numeric representation of the mood.

We used **production music** as ground truth data due to the size of the dataset (128,000 tracks) and the quality of the metadata.

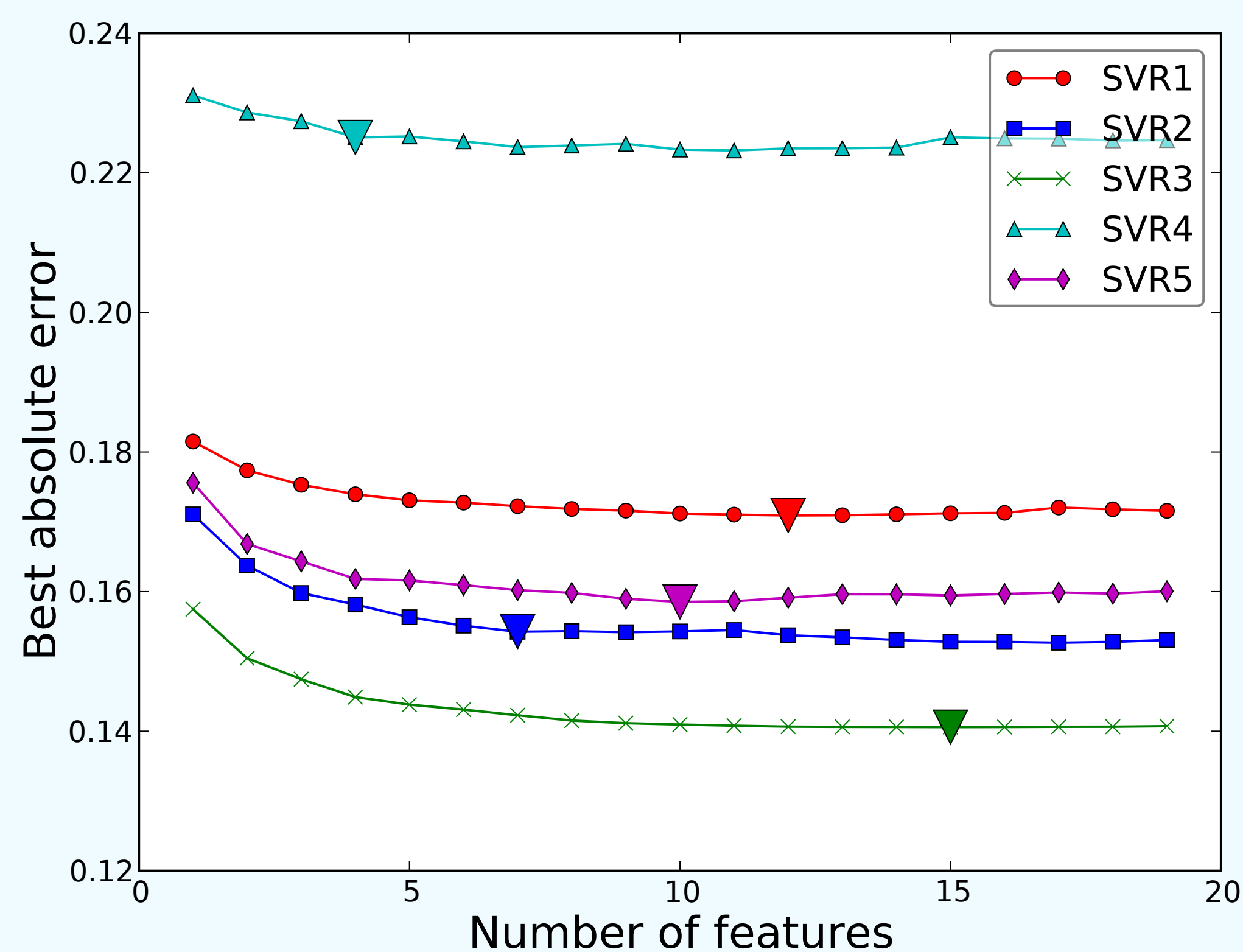
59 sets of features were extracted with a computing cluster using **Vamp plugins** and **Sonic Annotator**.

The mood representation was generated by taking the **structure** of the production music's mood metadata and applying **multi-dimensional scaling**. The optimum number of dimensions was found to be five through a subjective test [1].

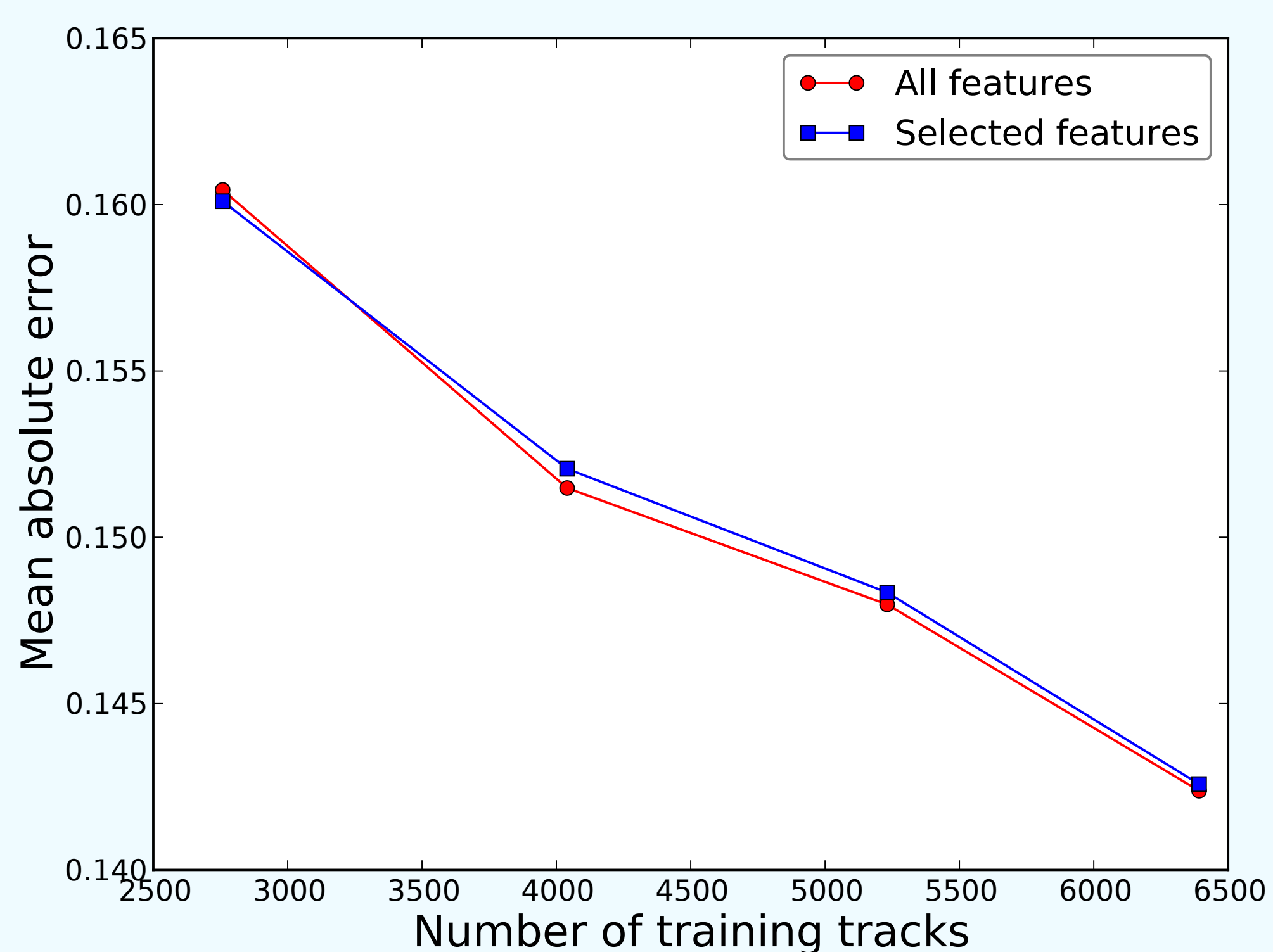


Production music mood keywords, arranged by frequency of co-occurrence.

A selected set of 1760 tracks was evaluated, using one-third for testing and 2-fold cross-validation. Every pair of features were tested before taking the top 12, testing those with every other feature and repeating.



The results show that **32 spectral, harmonic, rhythmic and temporal features** are needed for optimum performance, but as the error converges quickly, good performance can be achieved with much fewer.



The selected features match the performance of using all features and the error continues to drop after 6000 training tracks. Further work is needed to find the optimum number of training tracks.

Full results are available at bbcarp.org.uk/m4/aes53

[1] Mathieu Barthet, David Marston, Chris Baume, György Fazekas, and Mark Sandler. Design and Evaluation of Semantic Mood Models for Music Recommendation. In Proc. International Society for Music Information Retrieval Conference, 2013.