# Integration of mProphet chromatogram peak identification probability model into Skyline

Brendan MacLean[1], Don M. Marsh[1], Hannes Röst[2], Lucia Espona Pernas[2], Olga Schubert[2], George Rosenberger[2], Ruedi Aebersold[2], Michael J. MacCoss[1]

[1]University of Washington, Department of Genome Sciences, [2]Institute of Molecular Systems Biology, ETH Zurich

**Skyline**

## Overview:

The Skyline Targeted Proteomics Environment has distinguished itself as a reliable and useful tool for chromatography-based quantitative proteomics. From its initial focus on selected reaction monitoring (SRM) to its current support for full-scan methods including MS1 filtering, targeted MS/MS and data independent acquisition (DIA – including the approach popularized as SWATH) measurements of peptide abundance have been based on areas under peaks in mass spectrometric chromatograms. Until now, however, peak identification within Skyline has relied on a limited set of features without the ability to derive a statistical confidence metric like a false discovery rate (FDR). To correct this shortcoming, we have integrated the mProphet scoring features, semi-supervised learning (SSL) algorithm and derivation of peak q values into Skyline..

## Introduction:

Skyline now supports training an arbitrary number of chromatogram peak scoring models based on acquired data, using the mProphet SSL algorithm. These models consist of a set of coefficients, applied to available peak feature scores to derive an mProphet score, and a mean and standard deviation that define a normal distribution estimate of the null distribution for the mProphet scores. Different models may be trained for different experimental conditions. Once a model has been trained, it may be specified in Skyline and used for any number of experiments under the same conditions, without retraining under SSL. The null distribution parameters allow p value estimates, from which Skyline can estimate q values and FDR using the algorithm described by Storey, et al.

**Figure 1:** Peptide settings form showing a choice of trained mProphet scoring models as well as the legacy Skyline peak scoring model.

## Methods:

In order to achieve this, we have implemented the following:
- Decoy peptide and transition generation
- Decoy peptide and transition import
- A calculator architecture allowing two calculator types
  1. Summary – requiring only summary peak attributes
  2. Detailed – requiring chromatogram points
- Calculating and storing detail scores during import
- Calculators for the original 8 mProphet scores
  - ✓ Intensity[1]
  - ✓ Intensity correlation (dot-product) [1]
  - ✓ Coelution[2]
  - ✓ Shape[2]
  - ✓ Reference correlation (dot-product) [1]
  - ✓ Retention time deviation[1]
  - ✓ Reference coelution[2]
  - ✓ Reference shape[2]
- Calculators for several new scores
  - ✓ Coelution count[2]
  - ✓ Reference coelution count[2]
- Export mProphet features for testing with mProphet in R
- The mProphet SSL algorithm converted from R to C#
- Storey-Tibshirani q value calculation algorithm

In all cases, the resulting code was tested against data and results from original publications. Unit tests from OpenSWATH were used to validate some of the scores.

We have created two test Skyline documents:
1. The mProphet gold standard data set
2. A new document with decoy transitions generated in Skyline

Both have had SRM data imported and the new mProphet model editor used to train score coefficients for a composite linear scoring function.

**Figure 2a & b:** Original mProphet gold standard transition list from supplementary information that was imported into Skyline (left) and decoy generation form used to create decoy transitions in new document.

## Results:

**Figure 3a & b:** Target and decoy peptides, precursors, and transitions from the original mProphet gold method (left), where transitions were imported from the manuscript transition list, and a newly created method, where decoys were generated in Skyline (right).

**Figure 5:** Form for exporting peak groups and features to a CSV file for use with mProphet in R.

**Figure 5:** Nine runs imported normally into Skyline for the newly created document. Skyline plots show chromatograms with peaks picked and predicted retention times, matching MS/MS library spectrum, linear regression between measured retention times and iRT calculator scores and light and heavy precursor peak areas by run. The linear regression graph indicates that retention time prediction is unlikely to be a strong indicator of peak validity, since even with decoy peaks the R for the regression is 0.9986, probably because acquisition was scheduled in 4 minute windows around the predicted time.

**Figure 6:** Skyline scoring model training form, showing a fully mProphet model (top) with decoy and target histograms and the estimated decoy normal distribution. The trained linear model coefficients are shown (bottom-left) as well as the single score histograms for the bootstrap Intensity score (bottom-right)

**Figure 6:** Target-decoy histograms for the remaining 9 scores for the fully trained model shown in Figure 5, copied from the model training form.

**mProphet Gold Data** — **New Data**

**Figure 6:** Target-decoy and modeled null distribution graphs copied from the Skyline scoring model training form show improvement in separation of decoy and target peaks from the legacy Skyline peak score, after training with mProphet. Good separation persists, when a model trained on one data set is used on another. Null distribution estimates do not transfer well in this case.

## Conclusions:

- Good progress has been made on integrating the mProphet scoring model into Skyline in a way that makes important information for assessing model accuracy immediately available.
- Much work remains to be done to understand how much scoring separation will enhance peak picking in Skyline, and whether adding further scores can improve peak picking for non-SRM acquisition.
- Using calculated models to actually perform peak picking and q value estimation will require further implementation work in Skyline.
- Estimating accurate p and q values decoy distributions as a model for the null distribution may be an elusive goal.

**References:**
(1) Reiter L, Rinner O, et al. Nat. Methods. 2011/05; 8(5):430-5.
(2) Malmstrom L, Malmstrom J, et al. J. Proteome Res. 2012/03; 11(3):1644-53.