DART-ID Increases Proteome Coverage with Data-Driven Alignment of Retention Times

Albert Chen¹, Alexander Franks², Nikolai Slavov¹

1 Departments of Biology, Bioengineering, Northeastern University, Boston, MA 02115, USA **2** Department of Statistics and Applied Probability, UC Santa Barbara, CA 93106, USA

Summary

- Retention time (RT) is an informative feature that can be remarkably consistent across LC runs with the same sample and same experimental conditions.
- Data-driven alignments of peptide RTs across experiments create robust inferences of peptide RTs with RT distributions.
- Applying inferred and observed RTs within a principled Bayesian framework greatly increases the coverage of single cell proteomes.

Global RT Alignment Provides Robust RT Inferences Across Experiments



- Latent Reference RT (μ) inferred for each peptide sequence.
- Shifts between experiments modeled as monotonic transformation of μ .
- Optimization robust to large amounts of missing data

Global RT Alignment Results in Residuals < 1 min for 60 min LC Runs



Summary	statistics	for	residual	
-				

Summary statistics for residual RT (Δ RT)						
	Method	$ \Delta RT $	Median $ \Delta RT $			
ion	SSRCalc	3.77	3.19			
Predict	BioLCCC	3.44	2.91			
	ELUDE	2.51	2.02			
Alignment	iRT	0.693	0.537			
	MaxQuant	0.329	0.243			
	DART-ID	0.155	0.044			

- 46 LC-MS/MS runs, 60 min each
- Residual RT = Observed RT -Predicted RT
- DART-ID shows smallest average deviation of 0.044 min (2.6 seconds)

Residual RT varies by experiment



Bayesian Framework for Updating ID Confidence

$P(\text{ ID correct } | \text{ RT }) = \frac{P(\text{ RT } | \text{ ID correct }) \times P(\text{ ID correct })}{P(\text{ RT })}$

- P(ID correct | RT) Posterior probability that the PSM is assigned to the right sequence, given its RT. AKA "DART-ID" confidence.
- P(RT | ID correct) Conditional likelihood of the PSM's RT if its assigned sequence is correct. Estimated by evaluating inferred RT distribution of the sequence at the observed RT.
- P(RT) Marginal likelihood for observing the RT. Estimated as the sum of the probabilities that the PSM's assigned sequence is correct and incorrect.



ID Confidence Update Yields 50% More Peptides



Slavov Laboratory

Quantitative Biology



- Mean of inferred RT
- distribution is the transformed reference RT
- Inferred RT distribution
- variance experiment-specific
- Null distribution empirical
- distribution of all RTs
- Case 1 confidence upgraded Case 2 - confidence
- downgraded



DART-ID PSMs Give Consistent Protein Quantification





Conclusion

DART-ID takes advantage of reproducible retention times for peptide sequences within sets of LC runs to greatly increase the coverage of single cell proteomes. Global alignment method provides more robust estimates of RT. The more consistent experiment RTs are, the more powerful the added RT evidence is.

Acknowledgements

We thank members of the Slavov laboratory for discussions and feedback. This work was funded by the National Institute of Health to N.S. under Project Number 1DP2GM123497-01.

bioRxiv DOI: 10.1101/3991



Split data into disjoint sets

- **Spectra** PSMs filtered with PEP < 1%.
- **DART-ID** PSMs filtered with PEP > 1% and updated PEP < 1%. This set of new observations is *disjoint* from Spectra PSM.
- **Percolator** PSMs filtered with PEP > 1% and updated PEP < 1%. Also <u>disjoint</u> from Spectra PSMs.
- **Decoy** PSMs filtered with PEP < 1%, and randomly sampled to create a decoy protein.

Separation of the proteomes of 375 single cell samples each of 2 blood cancer cell lines, Jurkat and U-937 by Principal Component Analysis (PCA)