Using Machine Learning to Predict Recombinant Protein Expression

Abel Sousa^{1,5}, Stephanie K. Ashenden^{1,5}, Avid Afzal^{1,5}, Sergio Martinez Cuesta^{1,5}, Miguel Gancedo Rodrigo², William Lee³, Sandeep K. Talapatra^{2,4}, Dilrini De Silva¹, Rick Davies², Yinhai Wang¹, Ian Barrett¹, Aurelie Bornot^{1,6}, Lovisa Holmberg Schiavone^{3,6}

¹Data Sciences and Quantitative Biology, Discovery Sciences, AstraZeneca, Cambridge, UK, ²Discovery Biology Protein Sciences, Discovery Sciences, Cambridge, UK, ³Discovery Biology Protein Sciences, Discovery Biology Protein Sciences, Discov Sciences, Gothenburg, Sweden, ⁴Present address: Protein Science Group, Department of Protein and Cellular Sciences, GSK, Stevenage, UK, ⁵First co-authors, ⁶Corresponding authors: lovisa.holmberg.schiavone@astrazeneca.com

Abstract

The production of recombinant proteins is critical across several drug discovery stages. The process is costly and lengthy with a minimum of 6 weeks per construct to expression screen campaign, involving multiple steps. We are developing a machine learning platform that uses the primary sequence of proteins represented as physicochemical properties and structural features to support protein scientists by facilitating the design of protein constructs and highlighting sequences expressing at different yield classes. The model is coupled to an *in-silico* screening procedure that systematically designs and assesses thousands of constructs in a high-throughput manner. This method is currently being deployed in drug discovery projects and leads to the design of constructs expressing at a higher yield compared to those designed using human knowledge only. Here we share our initial results and plans to improve and develop the techniques through integrative teamwork and additional resources. We plan to do this by: (1) considering yield values instead of classes aided by GelClick, an automated gel image analysis tool, (2) incorporating deep learning features for sequence representation, and (3) leveraging external public domain datasets. Limited data to train the model is a key blocker so we are putting together a proof-of-concept and a pre-competitive consortium with academic and pharmaceutical industry partners to share data and models in collaboration with EMBL-EBI.

1 – The evolution of our predictive model

2017-2019

Exploring data for model building and implementation of expression screen module in AZ's lab management software.

Build initial model in an OŤ collaboration with an external partner. The model was put into production but gained little traction at the time.

ğuru supp													
												Recently Viewed	Sear
Expre 2022 2022 - Septemb	ssion screens AZ											_	
2022-09-0)6 - 0	coli											
ID: 640 🛔 Ali Al-B	ehadili 🖀 You and	200 more can edit									🏴 Flag Experime	nt ✔Sign More ▼	
Start date: Sep 06	2022 09: O Durati	on (days): Enter dura	tior									Calendar View	
🗣 tag this item													
Experimental of	letails											-	
E.coli cells expressi	on screen												
Assigned Scientist	Start Dat	0	End Date	В	ioELN Entry								
Stefan Blaho	26/08/2	022 🗖	06/09/2022	•	BioEln Number								
						add a section							
Table of constr E.coli cells expression	ructs on screen											•	
Row number	Gene	Plasmid		Plasmid ID	Mw (Da)	E.coli strain		Media		Co- expression	Induction	Antibiotic selection	
	IRF5				9986.72	BL21(DE3) Gold	× ×	ZYP-8012	× *	No × *	Autoinduction ×	* Kan × *	
1					8031.61	BL21(DE3) Gold	× *	ZYP-8012	× *	No × *	Autoinduction ×	* Kan × *	
2	IRF5									nic	1 Anna Anna Anna Anna Anna Anna Anna Ann	1 .	
2	IRF5 IRF5				9855.53	BL21(DE3) Gold	× *	ZYP-8012	× Ŧ	No × *	Autoinduction ×	* Kan × *	
1 2 3 4	IRF5 IRF5 IRF5				9855.53 7326.63	BL21(DE3) Gold BL21(DE3) Gold	× * × *	ZYP-8012 ZYP-8012	× * × *	No × •	Autoinduction ×	 Kan x * Kan x * 	

Improvements to the model based on lab data with **155** protein targets and more than 1100 constructs.

2020

Discovered that small sequence changes in constructs can result in large changes in production yield.

Implementation of **GelClick** – an Annotation of more data with the automated gel image analysis use of GelClick. tool.

Model piloted on 6 projects predicting and validating designed constructs.

2021

Initial exploration of deep learning embeddings by a postdoc (6 months).

MSc student analyses **public** data to be incorporated in the model and expands deep learning embeddings work.

2022

Initiating recruitment of a **postdoc** with EMBL-EBI and consolidation of a precompetitive **consortium**.

2 – Our current model

Calculate Input sequences features N-C-C - N-C-C - N-C-C Currently: AZ's lab data software Polar Physicochemical features Future: Secondary structures Protein Structure Initiativ

Output yield Predictive model category Prediction of yield category in an expression

host using a Random Forest model trained on lab data

- 1139 constructs, 155 proteins
- Assessment using 5-fold cross validation

Yield category	Precision (proportion of correct prediction)	Recall (proportion of true cases identified)
0 (0-1 mg/L)	0.48	0.63
1 (1-10 mg/L)	0.42	0.45
2 (10-20 mg/L)	0.40	0.04
3 (>20 mg/L)	0.37	0.19

3 – GelClick

A more precise yield estimation can help the predictive model to better learn the relationship between input sequence and yield. GelClick is a tool that accurately estimates the yield of recombinant proteins after expression screens and aims to select the best construct for production.

GelClick interface: The user selects the reference band (red) and the construct bands (green), enters parameters of the experiment and GelClick analyses the image components to estimate the yield of the constructs

Quantification of the estimation variability for the two methods: The disorder or scattering of









Global performance Random performance Metric 0.28 0.42 Accuracy

Accuracy better than random, OK at predicting lowly expressing constructs, there is room for improvement

estimations is lower when using the application. That can be seen in the orange curve that has a higher density around zero. A lower value of disordered is better since it means higher consistency.

4 – Project case study and summary of results

Disorder and flexibility¹



Our target is a transcription involved factor the in development of immune cells, B/T including cells and dendritic cells.

Its deregulation is involved in the development of different types of cancers.

Construct design

Protein constructs are designed *in* silico by systematically combining multiple protein regions; protein tags; and removing residues from the C/N-terminal domains.



Target regions considered: DBD, IAD and Full-length

Protein tag	Terminus
6His-GS-TEV-Target	Ν
6His-GS-ZZ-GS-TEV-Target	Ν
6His-SUMO-Target	Ν
GST-TEV-Target	Ν
М	Ν
GS-AVI	С
None	С

Up to 10 residues removed from the N/C-terminus

>3,500 constructs designed.

Predictions

Bioinformatics features of the in silico constructs are extracted using an inhouse pipeline.

ML model The can current provide expression predictions for 8 hosts.

			E.coli - bl21(de3)gold
Organism	Host		
E. coli	bl21(de3)gold	2000	
E. coli	bl21(de3)codonplus	1500	
E. coli	bl21(de3)rosettaplyss	nt	
E. coli	bl21(de3)star	<u></u> <u></u> <u></u> <u></u> <u></u> <u></u> <u></u> <u></u> <u></u> <u></u> <u></u> <u></u> <u></u> <u></u>	
E. coli	arcticexpress	500	
Insect	sf21		
Insect	sf9	0	
Mammal	expi293		Predicted yield

Yield class / Yield (mg/L): 0 / 0-1; 1 / 1-10; 2 / 10-20; 3 / >20

Ranking

The constructs are ranked using a desirability function that transforms and integrates multiple ranking features into a 0-1 scale, such that high values are desirable³.



predicted is NetSolP⁴ model.

been

Validation

are validated The constructs using expression screenings to confirm the predictions.



5 – Future work



References

1. Liu *et al.*, Molecules, 23(10), 2535 (2018) 2. Jumper et al., Nature, 596, 583–589 (2021) 3. Lazic S. E., PeerJ, 3:e1444 (2015) 4. Thumuluri *et al.*, Bioinformatics, 38(4), 941–946 (2022) 5. Detlefsen *et al.*, Nature Communications, 13, 1914 (2022)

Acknowledgements

We thank the members of the Data Sciences and Quantitative Biology (QuBi) and Discovery Biology departments in AstraZeneca for helpful inputs throughout the duration of the project and while creating this poster.

