

Bernhard Blank-
Landeshammer

Development and application of
proteogenomics methods to refine
genome annotations and detect
single amino acid variants

Development and application of proteogenomics methods to refine genome annotations and detect single amino acid variants

Zur Erlangung des akademischen Grades eines

Dr. rer. nat.

von der Fakultät Bio- und Chemieingenieurwesen
der Technischen Universität Dortmund
genehmigte Dissertation

vorgelegt von

Bernhard Blank-Landeshammer, MSc

aus

Kirchdorf a.d. Krems, Österreich

Tag der mündlichen Prüfung: 18.12.2020

1. Gutachter/-in: Prof. Dr. Sickmann
2. Gutachter/-in: Prof. Dr. Nett

Dortmund 2020

Berichte aus der Biologie

Bernhard Blank-Landeshammer

**Development and application of proteogenomics
methods to refine genome annotations and
detect single amino acid variants**

D 290 (Diss. Technische Universität Dortmund)

Shaker Verlag
Düren 2021

Bibliographic information published by the Deutsche Nationalbibliothek

The Deutsche Nationalbibliothek lists this publication in the Deutsche Nationalbibliografie; detailed bibliographic data are available in the Internet at <http://dnb.d-nb.de>.

Zugl.: Dortmund, Technische Univ., Diss., 2020

Copyright Shaker Verlag 2021

All rights reserved. No part of this publication may be reproduced, stored in a retrieval system, or transmitted, in any form or by any means, electronic, mechanical, photocopying, recording or otherwise, without the prior permission of the publishers.

Printed in Germany.

ISBN 978-3-8440-7854-1

ISSN 0945-0688

Shaker Verlag GmbH • Am Langen Graben 15a • 52353 Düren

Phone: 0049/2421/99011-0 • Telefax: 0049/2421/99011-9

Internet: www.shaker.de • e-mail: info@shaker.de

Abstract

Genome annotation procedures can benefit from information provided by proteomics in various ways. Detection of sample- or patient specific alterations of protein sequences in the form of single amino acid variants (SAAVs) require dedicated analysis and validation procedures. In this work, proteogenomics methods were established that enable proteomics analysis of organisms void of reference sequence databases and that can universally be used to refine genome annotations based on protein level information as well as for identification of SAAVs.

First, a *de novo* peptide sequencing method was established to improve database-free peptide identification. By integration of multiple algorithms, the reliability of the results could be improved and the number of identified sequences at stringent conditions was tripled compared to a reference analysis. The applicability of *de novo* peptide sequencing to differentiate species of foraminifera (unicellular protists) based on their proteome was assessed and a homology-based strategy was employed to quantify the proteomic response of the foraminifera *A. lobifera* to heat stress.

Next, the proteomic landscape of the freshwater snail *R. auricularia* was charted and a combination of *de novo* peptide sequencing and searches against the newly sequenced genome enabled an automated refinement of the annotation. This led to the additional identification of almost one thousand genes and major alterations of the underlying gene models. The genome annotation of the model fungus *S. macrospora* was subjected to similar refinements. Remarkably, *de novo* peptide sequencing enabled to our knowledge the first ever confirmation of A-to-I RNA editing in fungi on the protein level.

Finally, a proteogenomics approach was employed to analyze liver metastases retrieved from patients suffering from colorectal cancer. This enabled the identification of multiple SAAVs such as the prognostic KRAS G12V mutation. These discoveries were further validated and quantified by a targeted MS assay using of stable-isotope labeled standard (SIS) peptides.

Abstract (in German)

Genomannotation stellt in der Form von Protein Sequenzdatenbanken die Basis für LC-MS basierte Proteomanalysen zur Verfügung. Proteomische Daten können umgekehrt aber auch dazu beitragen, diese Annotationsvorgänge zu verbessern, indem Genmodelle verfeinert und Kodierungspotentiale besser bewertet werden können. In dieser Arbeit wurden proteogenomische Methoden etabliert um Proteomanalysen von Organismen mit fehlender Referenzdatenbank durchzuführen, die Annotation von Genomen mithilfe von Proteomdaten zu verbessern und Einzel-Aminosäureaustausche (SAAVs) sicher zu identifizieren.

Zuerst wurde eine kombinatorische *de novo* Peptidsequenziermethode etabliert, um eine Datenbank-freie Identifikation von Peptiden zu ermöglichen. Durch die Integration mehrerer Algorithmen konnte die Zuverlässigkeit der Ergebnisse verbessert werden und die Anzahl der identifizierten Peptide bei relativ stringenten Kriterien (d.h. 5% FDR) konnte im Vergleich zum besten Einzel-Algorithmus verdreifacht werden. Die Anwendbarkeit dieser *de novo* Sequenziermethode wurde anhand eines Speziesvergleichs von Foraminiferen (einzellige Protisten) getestet und mittels einer Homologie-basierten Suchstrategie wurde die Proteomantwort der Foraminiferen-Spezies *A. lobifera* auf Hitzestress quantifiziert.

Ferner wurde das Proteom der Frischwasserschnecke *R. auricularia* charakterisiert und eine Kombination aus *de novo* Sequenzierung und Datenbanksuchen gegen das translatierte Genom dazu genutzt, diese Informationen automatisiert in die Annotation einfließen zu lassen. Dies erlaubte die Identifikation von annähernd 1,000 Genen und eine beträchtliche Abänderung der zugrundeliegenden Genmodelle. Dieser Arbeitsablauf wurde weiter dazu genutzt, die Genomannotation des als Modelorganismus genutzten Pilzes *S. macrospora* zu verbessern. Hier ermöglichte der *de novo* Sequenzier-Ansatz den unseres Wissens nach ersten Nachweis von A-zu-I RNA Editierung in Pilzen auf Proteinebene.

Schließlich wurde ein proteogenomischer Ansatz dazu verwendet, Gewebeproben von Lebermetastasen zu analysieren, die Patienten mit kolorektalem Karzinom entnommen wurden. Dabei wurden mehrere SAAVs auf Proteinebene nachgewiesen, unter anderem die KRAS G12V Mutation. Eine zielgerichtete MS Untersuchung mit stabil isotopenmarkierten Standardpeptiden wurde etabliert um den Mutationsgrad zu quantifizieren.

Table of contents

Abstract	2
Abstract (in German)	3
Table of contents	4
Abbreviations	1
Proteinogenic amino acids	3
List of publications	5
List of conference contributions	6
Talks	6
Poster presentations	6
Declaration of pre-published contents	7
1. Introduction	8
1.1. Genomics.....	8
1.1.1. <i>Eukaryotic genome and gene structure</i>	8
1.1.2. <i>Genome sequencing and annotation</i>	14
1.1.3. <i>Computational gene prediction</i>	16
1.2. Proteome analysis.....	17
1.2.1. <i>LC-MS based proteomics</i>	18
1.2.2. <i>Electrospray ionization</i>	22
1.2.3. <i>Mass analyzers</i>	23
1.2.4. <i>Tandem mass spectrometry for proteomics</i>	26
1.2.5. <i>Data acquisition strategies</i>	28
1.2.6. <i>Quantitative proteomics</i>	29
1.2.7. <i>Proteomics MS data analysis</i>	32
1.3. Proteogenomics	36
1.3.1. <i>Environmental proteogenomics and metaproteomics</i>	37
1.3.2. <i>Proteogenomics for genome annotation</i>	39
1.3.3. <i>Proteogenomics in Precision medicine</i>	41
2. Aim	45
3. Material and method	46
3.1. Materials	46
3.1.1. <i>Chemicals</i>	46
3.1.2. <i>Instruments and chromatography columns</i>	47
3.1.3. <i>Disposable Consumables</i>	48

3.1.4. <i>Data analysis software</i>	49
3.2. <i>Methods</i>	50
3.2.1. <i>Samples for proteomics analyses</i>	51
3.2.2. <i>Lysis, Homogenization and protein extraction</i>	52
3.2.3. <i>Determination of protein concentration</i>	53
3.2.4. <i>Carbamidomethylation</i>	53
3.2.5. <i>Filter-aided sample preparation</i>	54
3.2.6. <i>Ethanol precipitation and in-solution digestion</i>	54
3.2.7. <i>Evaluation of digestion efficiency</i>	55
3.2.8. <i>Desalting of proteolytic digests</i>	56
3.2.9. <i>High-pH RP offline fractionation</i>	56
3.2.10. <i>Phosphopeptide enrichment</i>	57
3.2.11. <i>Synthesis and purification of synthetic (SIS-) peptides</i>	58
3.2.12. <i>Nano LC-ESI-MS/MS analysis</i>	59
3.2.13. <i>Raw data analysis</i>	60
3.2.14. <i>Data processing, integration and statistical evaluation</i>	69
4. <i>Results</i>	75
4.1. <i>De novo peptide sequencing</i>	75
4.1.1. <i>Optimization of MS parameters</i>	75
4.1.2. <i>Evaluation of de novo peptide sequencing algorithms</i>	76
4.1.3. <i>Validation of de novo generated peptides sequences with synthetic peptides</i>	80
4.2. <i>Database-independent characterization of Foraminifera sp.</i>	81
4.2.1. <i>Homology-based characterization of the proteome-response to thermal-stress in A. gibbosa</i>	85
4.3. <i>Genome annotation of R. auricularia</i>	90
4.3.1. <i>Tripartite data analysis workflow</i>	91
4.3.2. <i>Automated refinement of the R. auricularia genome annotation</i>	94
4.4. <i>Refinement of the S. macrospora genome annotation by Proteogenomics</i>	97
4.4.1. <i>Identification of novel genes</i>	100
4.4.2. <i>Functional Co-Expression analysis</i>	103
4.4.3. <i>Alternative splice variants</i>	104
4.4.4. <i>Detection of SAAVs induced by RNA Editing</i>	106
4.4.5. <i>Absolute quantification of Stop-Loss editing variants using SIL peptides</i>	110
4.5. <i>Global proteome analysis and detection of SAAVs in Colorectal Cancer tissue</i>	111
4.5.1. <i>Validation of SAAV identifications by PRM</i>	115

5.	Discussion and Conclusion.....	121
5.1.	<i>De novo</i> peptide sequencing – a viable addition to database search approaches.....	121
5.2.	Proteomics analysis of Foraminifera.....	124
5.3.	<i>Radix auricularia</i> – Automated proteome-driven genome annotation.....	125
5.4.	Refinement of <i>S. macrospora</i> genome – new tricks for old genomes.....	128
5.5.	KRAS p.G12V mutation in liver metastases of colorectal cancer tissue	131
6.	References	135
7.	Acknowledgements	152
8.	Supplement.....	153
8.1.	LC-MS parameters.....	153
8.2.	Overview of <i>de novo</i> validation datasets.....	155
8.3.	Foraminifera species comparison	157
8.4.	<i>De novo</i> -aided genome annotation of <i>R. auricularia</i>	158
8.5.	<i>S. macrospora</i> RNA-Editing analysis.....	159
8.6.	Global analysis of CRC liver metastases	162
8.7.	Absolute quantification of SAAVs in CRC liver metastases	165
8.8.	List of supplementary data files	167