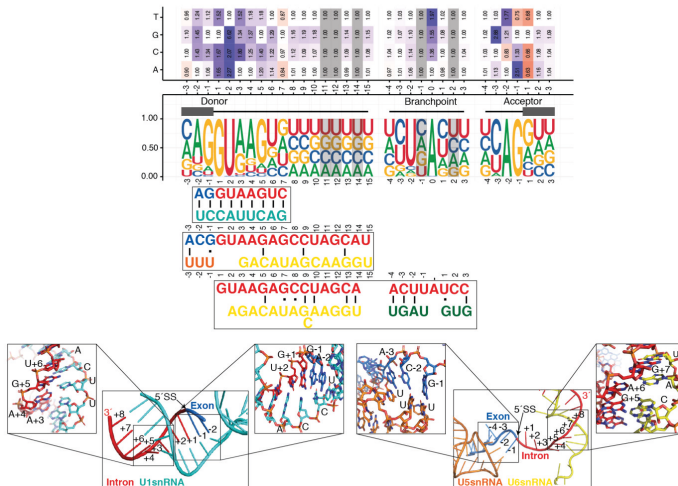


Leonhard Wachutka

Global Donor and Acceptor Splicing Site Kinetics in Human Cells





Fakultät für Informatik
Technische Universität München



Global Donor and Acceptor Splicing Site Kinetics in Human Cells

Leonhard Konrad Friedrich Wachutka

Vollständiger Abdruck der von der Fakultät für Informatik der Technischen Universität München zur Erlangung des akademischen Grades eines

Doktor der Naturwissenschaften (Dr. rer. nat.)

genehmigten Dissertation.

Vorsitzender und Prüfer der Dissertation:

Prof. Dr. Björn Menze
Technische Universität München

Prüfende der Dissertation:

1. Prof. Dr. Julien Gagneur
Technische Universität München
2. Prof. Dr. Caroline Friedel
Ludwig-Maximilians-Universität München

Die Dissertation wurde am 27.05.2020 bei der Technischen Universität München eingereicht und durch die Fakultät für Informatik 04.08.2020 angenommen.

Selected Topics of Electronics and Micromechatronics
Ausgewählte Probleme der Elektronik und Mikromechatronik

Volume 52

Leonhard Konrad Friedrich Wachutka

**Global Donor and Acceptor Splicing Site Kinetics
in Human Cells**

Shaker Verlag
Düren 2020

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.: München, Techn. Univ., Diss., 2020

Copyright Shaker Verlag 2020

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 978-3-8440-7710-0

ISSN 1618-7539

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

I. Acknowledgments

So finally this thesis comes to an end. If someone had told me five years ago that this ever would happen, I would not have believed any word. I would like to thank the countless people who made this possible. The people, who gave me hope, when I had no more. The people, who made impossible things possible. The people, who supported me to do it 'my way', although I know it was not the easiest. And all the people I met during the past five years, who inspired me never to give up. I want to say thank you to every single one of these many people, because nothing would have worked out without you.

In addition, I want to say special thanks to Julien Gagneur, my supervisor, who didn't let the faith go in me. He has been a constant source of new ideas and the best boss one can think of. And I thank Livia Caizzi: without her help and her endless patience during the experiments nothing of this would have become possible. I thank Patrick Cramer for all the useful input and support. I thank Carina Demel, who inspired me during my master thesis and who helped me with bioinformatics. I thank Jun Cheng, who jump started the bioinformatics part of this thesis and who always had good ideas in his mind. I thank Christian Mertes and Žiga Avsec, for countless hours of fruitful discussions and their continuous help. And I thank all people of the past and present Gagneurlab. You folks are great!

I want to express my special thanks to my whole family for their unbelievable support, my father who opened up my mind for natural science, my mother who always grounded me, and my lovely brother, who makes my life happier every day. And I want to express my particular thanks to Anna, Tins and Lexi for their endless love. Without love this live wouldn't be worth living.

II. Abstract

Extensive *in vitro* studies have strongly advanced our understanding of the splicing process, but the kinetics of the successive catalytic steps of splicing *in vivo* remains far less understood. Most existing studies fall short of a sound theoretical interpretation of measured data because of the entangled nature of RNA synthesis, splicing and degradation in standard RNA-Seq protocols.

TT-seq is a labeled RNA sequencing protocol combined with an early RNA fragmentation step that makes it possible to quantify the human RNA metabolism at the resolution of individual phosphodiester bonds. Using Transient Transcriptome sequencing (TT-seq) enabled, for the first time, the genome-wide quantification of the half-lives of individual donor site and acceptor site bonds as well as the rate of formation of splice junctions.

A careful analysis of the kinetic times provided the basis for a systematic search for genetic predictors of fast and slow splicing processes. It turned out that the donor site cleavage time is limited by polymerase elongation and allowed for estimating the interaction effects between several spliceosome components and the RNA at single nucleotide resolution. These findings agree very well with published structural data.

During this study it was necessary to set up solid mathematical models of single bond splicing kinetics to overcome the experimental ambiguities caused by alternative splicing, which are common to massive parallel sequencing methods. This approach, for the first time, allows for introducing and measuring the splicing yield, which is the proportion of precursor RNA successfully converted into spliced RNA.

With a view to facilitating similar analyses for other researchers, the software package rCube has been developed, which opens an easy access to the analysis of TT-seq data to estimate RNA kinetics. It is not only applicable to infer splicing kinetics, but also allows for the high precision measurement of RNA turnover rates.

In summary, this thesis provides conceptual advances in the understanding of RNA kinetics, demonstrates the power of the approach by finding de-novo known and unknown genetic determinants of splicing, and offers a rich source of data for further analysis.

III. Publications

Global donor and acceptor splicing site kinetics in human cells

Leonhard Wachutka^{1†}, Livia Caizzi^{2†}, Julien Gagneur¹, Patrick Cramer²

¹Department of Informatics, Technical University of Munich, Garching, Germany;

²Department of Molecular Biology, Max-Planck-Institute for Biophysical Chemistry, Göttingen, Germany

†These authors contributed equally to this work

(2019) eLife. DOI: 10.7554/eLife.45056. (Wachutka *et al.*, 2019)

Author contributions: **Leonhard Wachutka**, conceptualization, data curation, software, formal analysis, validation, investigation, visualization, methodology, writing - original draft, writing - review and editing, designed and carried out the bioinformatics analysis; Livia Caizzi, conceptualization, data curation, validation, investigation, visualization, methodology, writing - original draft, writing - review and editing, optimized and carried out TT-seq experiments, contributed to the design of the bioinformatics analysis, and used molecular modeling to interpret results; Julien Gagneur, conceptualization, resources, software, supervision, funding acquisition, investigation, visualization, methodology, writing - original draft, project administration, Writing - review and editing; Patrick Cramer, conceptualization, resources, supervision, funding acquisition, investigation, visualization, methodology, writing - original draft, project administration, writing - review and editing author

Measures of RNA metabolism rates: Toward a definition at the level of single bonds

Leonhard Wachutka and Julien Gagneur

Department of Informatics, Technical University of Munich, Garching, Germany

(2017) Transcription. DOI: 10.1080/21541264.2016.1257972. (Wachutka and Gagneur, 2017)

Author contributions: **Leonhard Wachutka**, conceptualization, visualization, writing - original draft, writing - review and editing; Julien Gagneur, conceptualization, resources, supervision, funding acquisition, writing - original draft, project administration, writing - review and editing

Transient transcriptome sequencing: computational pipeline to quantify genome-wide RNA kinetic parameters and transcriptional enhancer activity

Gabriel Villamil^{1†}, **Leonhard Wachutka**^{2†}, Patrick Cramer¹, Julien Gagneur², Björn Schwalb^{1†}

¹Department of Molecular Biology, Max-Planck-Institute for Biophysical Chemistry, Göttingen, Germany

²Department of Informatics, Technical University of Munich, Garching, Germany;

†These authors contributed equally to this work

(2019) bioRxiv. DOI: 10.1101/659912. (Villamil *et al.*, 2019)

Author contributions: TT-seq protocol description, writing: Gabriel Villamil, Patrick Cramer, Björn Schwalb. rCube analysis description, writing: Leonhard Wachutka, Julien Gagneur

OCR-Stats: Robust estimation and statistical testing of mitochondrial respiration activities using Seahorse XF Analyzer

Vicente A. Yépez^{1,2}, Laura S. Kremer^{3,4}, Arcangela Iuso^{3,4}, Mirjana Gusic^{3,4}, Robert Kopajtich^{3,4}, Eliska Koňáriková^{3,4}, Agnieszka Nadel^{3,4}, **Leonhard Wachutka**¹, Holger Prokisch^{3,4}, Julien Gagneur^{1,2}

¹ Department of Informatics, Technical University of Munich, Garching, Germany,

² Quantitative Biosciences Munich, Gene Center, Department of Biochemistry, Ludwig-Maximilians Universität München, Munich, Germany,

³ Institute of Human Genetics, Helmholtz Zentrum München, Neuherberg, Germany, ⁴ Institute of Human Genetics, Klinikum Rechts der Isar, Technical University of Munich, Munich, Germany

(2018) PLoS ONE. DOI: 10.1371/journal.pone.0199938 (Yépez *et al.*, 2018)

Author Contributions:

Conceptualization: Vicente A. Yépez, Laura S. Kremer, Arcangela Iuso, Mirjana Gusic, Robert Kopajtich, Eliska Koňáriková, Agnieszka Nadel, Holger Prokisch, Julien Gagneur.

Data curation: Laura S. Kremer, Arcangela Iuso, Mirjana Gusic, Robert Kopajtich, Eliska Koňáriková, Agnieszka Nadel.

Formal analysis: Vicente A. Yépez, Holger Prokisch, Julien Gagneur.

Investigation: Vicente A. Yépez, Julien Gagneur. Software: Vicente A. Yépez, **Leonhard Wachutka**. Supervision: Holger Prokisch, Julien Gagneur.

Visualization: Vicente A. Yépez, Julien Gagneur. Writing – original draft: Vicente A. Yépez, Holger Prokisch, Julien Gagneur.

Writing – review & editing: Vicente A. Yépez, Laura S. Kremer, Arcangela Iuso, Mirjana Gusic, Robert Kopajtich, Eliska Koňáříková, Agnieszka Nadel, **Leonhard Wachutka**, Holger Prokisch, Julien Gagneur.

IV. Content

I. Acknowledgments	1
II. Abstract	3
III. Publications	5
IV. Content.....	9
1 Introduction.....	13
1.1 Gene expression.....	13
1.2 Spliceosome consists of many subunits	14
1.3 Splicing mechanism in detail.....	17
1.4 RNA splicing kinetics	18
1.5 Experimental background	20
1.5.1 RNA-seq for quantifying the transcriptome expression	20
1.5.2 4sU-seq for isolating nascent RNA	21
1.5.3 TT-seq for uniform mapping of the transient transcriptome.....	23
1.6 Aims and scope of this thesis	24
1.6.1 Define novel splicing metric.....	24
1.6.2 Modeling and measurement of RNA splicing kinetics.....	25
1.6.3 Finding the genetic determinants of splicing kinetics	25
1.6.4 Development of a software package for other researchers	26
2 Definition and Estimation of Splicing Quantities at the Single Bond Level	27
2.1 Splicing quantities at the single bond level	27
2.1.1 Definition of RNA metabolism at the level of single bonds	27
2.1.2 Steady-state splicing parametrization.....	29
2.1.3 Splicing yield	30
2.1.4 Splicing quantities found in literature.....	31
2.2 Modeling and assessment of genome-wide splicing kinetics	32

2.2.1	Advantage of TT-seq over 4sU-seq	32
2.2.2	TT-seq time series experiment	35
2.2.3	Kinetic RNA splicing models	38
2.3	Estimation of sample normalization factors and cross-contamination	50
2.4	Kinetic rate modeling and estimation	52
2.5	Estimation of the relative uncertainty for the kinetic parameters	55
3	Predictors of Splicing Quantities	57
3.1.1	Human mRNA splicing takes median of 7 min	57
3.1.2	Intron length constrains splicing times	59
3.1.3	Donor site kinetics depend on U1 and U5 snRNA interactions	62
3.1.4	Structural determinants of acceptor site kinetics	66
3.1.5	Splicing yield	68
4	The Software rCube	73
4.1	Architectural overview	73
4.2	Implementation	74
4.3	Usage of rCube	76
4.3.1	Overview	76
4.3.2	Installation and loading of the package	77
4.3.3	Defining genomic features	77
4.3.4	The rCubeExperiment class	78
4.3.5	Counting reads	79
4.3.6	Sample normalization and cross-contamination	79
4.3.7	Estimating the dispersion	80
4.3.8	Fitting rates	80
5	Conclusion and Discussion	83
5.1	Conceptual contributions	83
5.2	Biological insights	84

5.3	Conceptual improvements for motif discovery and validation	85
5.4	Software.....	86
5.5	Limitations.....	87
6	Appendix	91
6.1	Cell culture.....	91
6.2	TT-seq time series	91
6.3	Read alignment and counting	92
6.4	Determination of the major isoforms	93
6.5	Branchpoint identification.....	93
6.6	Estimation of single nucleotide effects.....	93
6.7	Comparison of 4sU-seq and TT-seq.....	94
7	List of Figures.....	95
8	List of tables	97
9	Bibliography.....	99

Parts of this thesis have already been published. The respective publications and the contributions of other co-authors are clearly indicated at the beginning of each chapter. The first person plural form 'we' is used throughout the thesis to avoid unnecessary switching between forms 'I' and 'we'.