

Dissertation zur Erlangung des Doktorgrades eingereicht am  
Fachbereich Mathematik und Informatik der Freien Universität Berlin

# The COMMON DENOMINATOR PROCEDURE

A Novel Approach to Gene Expression Data Mining for  
Identification of Phenotype-Specific Genes

René Korn

2006

Betreuer:

PD Dr. Steffen Schulze-Kremer

Dr. Sascha Röhrig

Dr. Ulrich Brinkmann

Gutachter:

PD Dr. Steffen Schulze-Kremer

Prof. Dr. Peter Buckel

Tag der Disputation:

28. April 2006

*Dedicated to my beloved father Klaus Korn,  
who had to suffer from cancer for too many years.*

# Abstract

This thesis addresses the gap between the amount of on-hand expression data and the availability of information related to the function of those genes. To this end, a data mining procedure for the identification of genes that are associated with pre-defined phenotypes and/or molecular pathways was established. Based on the observation that pathway/phenotype associated genes are frequently expressed in same or nearby places and at identical or similar time points, an approach termed COMMON DENOMINATOR PROCEDURE (CDP) was devised. One unique feature of this novel approach is that the specificity and probability to identify desired phenotype/pathway-associated factors increases the more diverse the input data are. Three different approaches are discussed and compared: (i) a BASIC CDP, (ii) a GENETIC ALGORITHM BASED CDP and (iii) an INDICATOR GENES BASED CDP. To show the feasibility of these approaches, the CGAP Expression Data combined with a defined set of angiogenic factors was used to identify additional and novel angiogenesis-associated genes. A multitude of these additional genes were known to be associated with angiogenesis according to published data, verifying our approach. Application of a high throughput functional genomics platform (XantoScreen™) provided further experimental evidence for association of candidate genes with angiogenesis.

# Acknowledgement

First of all, I would like to thank PD Dr. Steffen Schulze-Kremer for being my supervisor and for supporting this work with his beneficial ideas and constructive criticism. I would also like to kindly thank Prof. Dr. Peter Buckel for taking over the part of the second referee. I would like to thank Dr. Ulrich Brinkmann for the original impulse to this work and his commitment.

I also appreciate all my friends and colleagues from Xantos, for their helpful feedback and their friendship. I want to thank Dr. Irene Boche, Stefan Hess, Pii Lämmlein, Dr. Dieter Link, Dr. Stephan Reschauer and Dr. Jürgen Zitzler for their friendship, a nice time and exciting badminton matches. Special thanks go to Dr. Sascha Röhrig, my patient superior, mentor and friend, as well as my closest colleagues from the Bioinformatics and IT department, Dr. Bettina Ehring, Dr. Beate Gawin, Fabian Weiß, Alexander Felber, Dr. Björn Kesper, Dr. Reinhold Köckerbauer and Dr. Alexander Spychaj. I am particularly grateful for the support from the high throughput screening team, Johannes Görl, Dr. Kerstin König-Hoffmann, Dr. Rolf Schäfer and Michael Kazinski.

I want to demonstrate my deepest respect to Xantos Biomedicine AG, especially to Stephan Wehselau and again to Prof. Dr. Peter Buckel, for their investment in my future and their commitment during exciting, interesting and turbulent times.

I am also thankful for the possibility to submit my work to the Department of Mathematics and Computer Science at the Free University of Berlin.

Last but not least I want to thank my family, especially my partner in life Janine Christ, my mother Brigitte Korn, my sister Carolin Korn and my godfather Rainer Spitzenpfeil as well as all my friends, especially Roman Egle and Daniel Godau, for encouragement, for proofreading, for putting up with me all the time, and for just being there.

# Contents

Abstract . . . . .	i
Acknowledgement . . . . .	ii
Contents . . . . .	iii
List of Figures . . . . .	v
List of Tables . . . . .	vi
Abbreviations . . . . .	vii
<b>1 Introduction</b>	<b>1</b>
1.1 Knowledge Discovery in Databases . . . . .	1
1.2 Biological Background . . . . .	4
1.2.1 Central Dogma of Molecular Biology . . . . .	4
1.2.2 Gene Expression . . . . .	7
1.3 Biological Challenge . . . . .	7
1.3.1 Common Denominator Concept . . . . .	8
1.3.2 Phenotype Angiogenesis . . . . .	9
1.4 Outline . . . . .	10
<b>2 System and Methods</b>	<b>12</b>
2.1 Infrastructure . . . . .	12
2.2 Data Sources . . . . .	13
2.2.1 Adaptation to CGAP Expression Data . . . . .	13
2.2.2 Definition of the INDICATORGENESET . . . . .	14
2.2.3 Definition of the ANGIOTESTGROUP . . . . .	14
2.3 COMMON DENOMINATOR PROCEDURE (CDP) . . . . .	15
2.3.1 Generation of the LIBRARYPROFILE . . . . .	17
2.3.2 Calculation of the GENESCORE . . . . .	20
2.3.3 Control Profiles . . . . .	22
2.3.4 Definition of ANGIOPROFILES . . . . .	23
2.3.5 Selection of Candidate Genes . . . . .	23
2.4 XantoScreen™ . . . . .	25
2.4.1 HUVEC Proliferation High Throughput Screening Assay . . . . .	27

<b>3</b>	<b>Results</b>	<b>29</b>
3.1	BASIC CDP . . . . .	30
3.1.1	Definition of Input Data . . . . .	30
3.1.2	Determination of the LIBRARYPROFILE . . . . .	32
3.1.3	Determination of the GENESCORE . . . . .	32
3.1.4	Selection of Candidate Genes . . . . .	33
3.1.5	Procedure Control and Validation . . . . .	33
3.2	GENETIC ALGORITHM BASED CDP . . . . .	37
3.2.1	Definition of Input Data . . . . .	37
3.2.2	Determination of the LIBRARYPROFILE . . . . .	39
3.2.3	Determination of the GENESCORE . . . . .	39
3.2.4	Selection of Candidate Genes . . . . .	39
3.2.5	Procedure Control and Validation . . . . .	40
3.3	INDICATOR GENES BASED CDP . . . . .	43
3.3.1	Definition of Input Data . . . . .	43
3.3.2	Determination of the LIBRARYPROFILE . . . . .	45
3.3.3	Determination of the GENESCORE . . . . .	45
3.3.4	Selection of ANGIOPROFILES . . . . .	46
3.3.5	Selection of Candidate Genes . . . . .	46
3.3.6	Procedure Control and Validation . . . . .	48
3.4	Summary . . . . .	53
3.4.1	Internal Procedure Control . . . . .	54
3.4.2	Procedure Validation - Experimental . . . . .	55
3.4.3	Procedure Validation - Literature . . . . .	56
<b>4</b>	<b>Discussion</b>	<b>57</b>
4.1	Comparison of the Procedures . . . . .	58
4.2	Comparison to Established Procedures . . . . .	60
4.3	Extensibility . . . . .	61
4.4	Future Perspective . . . . .	63
	<b>References</b>	<b>66</b>
<b>A</b>	<b>Data Sources</b>	<b>75</b>
<b>B</b>	<b>Implementation</b>	<b>78</b>
<b>C</b>	<b>Anhang gemäß Promotionsordnung</b>	<b>80</b>
C.1	Erklärung . . . . .	80
C.2	Lebenslauf . . . . .	81
C.3	Zusammenfassung . . . . .	82

# List of Figures

1.1	Different Steps of Knowledge Discovery in Databases . . . . .	2
1.2	Central Dogma of Molecular Biology . . . . .	5
1.3	Different Steps of Protein Synthesis . . . . .	6
2.1	Key Steps of the IGCDP . . . . .	16
2.2	Example of a Genetic Algorithm . . . . .	19
2.3	Schematic Candidate Gene Selection . . . . .	26
3.1	Flowchart of the BASIC CDP. . . . .	31
3.2	Flowchart of the GENETIC ALGORITHM BASED CDP. . . . .	38
3.3	Flowchart of the INDICATOR GENES BASED CDP. . . . .	44
3.4	Validation of the INDICATOR GENES BASED CDP. . . . .	51
4.1	Distribution of UniGene Clusters within CGAP Expression Data . . . . .	62
B.1	Entity Relationship Diagram . . . . .	78
B.2	UML Class Diagram . . . . .	79

# List of Tables

2.1	Number of ANGIOTESTGROUP Genes in the Selected CGAP Libraries . . .	18
3.1	Top Candidate Genes of the BASIC CDP . . . . .	34
3.2	Expectation of ANGIOTESTGROUP Genes and Screen Hits . . . . .	35
3.3	Number of Candidate Genes of the BASIC CDP . . . . .	36
3.4	Top Candidate Genes of the GENETIC ALGORITHM BASED CDP . . . . .	41
3.5	Number of Candidate Genes of the GENETIC ALGORITHM BASED CDP . .	42
3.6	Detailed Information for the Selected ANGIOPROFILES . . . . .	47
3.7	Top Candidate Genes of the INDICATOR GENES BASED CDP . . . . .	49
3.8	Number of Candidate Genes of the INDICATOR GENES BASED CDP . . . .	50
3.9	Known Modulators of Angiogenesis for the INDICATOR GENES BASED CDP	53
4.1	Comparison of the COMMON DENOMINATOR PROCEDURES . . . . .	59
A.1	CGAP Libraries Selected as Data Source . . . . .	76
A.2	Tissue Distribution of the Selected CGAP Libraries . . . . .	77
A.3	Protocol Distribution of the Selected CGAP Libraries . . . . .	77
A.4	Histology Distribution of the Selected CGAP Libraries . . . . .	77



# Abbreviations

BCDP	BASIC COMMON DENOMINATOR PROCEDURE
cDNA	Complementary DNA
CDP	COMMON DENOMINATOR PROCEDURE
CGAP	NCI's Cancer Genome Anatomy Project
CPU	Central Processing Unit
DDIT4	HIF-1 Responsive RTP801
DMEM	Dulbecco's Modified Eagle Medium
DNA	Deoxyribonucleic Acid
ECBM	Endothelial Cell Basal Medium, Promocell®
ECGF	Endothelial Cell Growth Factor 1
ECGM	Endothelial Cell Growth Medium, Promocell®
EPAS1	Endothelial PAS Domain Protein 1
ERD	Entity Relationship Diagram
EST	Expressed Sequence Tag
FDA	U.S. Food and Drug Administration
FCS	Foetal Calf Serum, Invitrogen®
GA	Genetic Algorithm
GB	Gigabyte
GHz	Gigahertz
GACDP	GENETIC ALGORITHM BASED COMMON DENOMINATOR PROCEDURE
GO	Gene Ontology
GRIF	Gene References Into Function
HEK293	Human Embryonal Kidney Cells
HIF1A	Hypoxia-inducible Factor 1, Alpha Subunit
HUVEC	Human Umbilical Cord Vein Endothelial Cells
IGCDP	INDICATOR GENES BASED COMMON DENOMINATOR PROCEDURE
IGFR1	Insulin-like Growth Factor 1 Receptor
JGAP	Java Genetics Algorithms Package
MB	Megabyte
mRNA	Messenger RNA
NCBI	U.S. National Center for Biotechnology Information
NCI	U.S. National Cancer Institute
NIH	U.S. National Health Institute
NHGRI	U.S. National Human Genome Research Institute
PDGF	Platelet-derived Growth Factor

RAM	Random Access Memory
RNA	Ribonucleic Acid
SAGE	Serial Analysis of Gene Expression
VEGF	Vascular Endothelial Growth Factor
UniProt	Universal Protein Resource
RefSeq	Protein Sequence Databases (Reference Sequence)
GenBank	Nucleotide Sequence Database
EMBL	European Molecular Biology Laboratory
UML	Unified Modeling Language