

Konrad Büssow

**Arrayed cDNA expression libraries for antibody
screening and systematic analysis of gene products**

Dissertation

**eingereicht im Fachbereich Chemie
der Freien Universität Berlin**

1998

1. Gutachter: Prof. Dr. H. Lehrach

2. Gutachter: Prof. Dr. V. Erdmann

Die vorliegende Arbeit wurde in der Zeit von Mai 1995 bis November 1998 in der Abteilung Prof. Dr. H. Lehrach am Max-Planck-Institut für Molekulare Genetik in Berlin durchgeführt. Ich möchte mich bei Herrn Prof. Lehrach bedanken, der die Entstehung dieser Arbeit in seiner Abteilung in großzügiger Weise unterstützt hat und mir mit Ratschlägen zur Seite stand. Besonderer Dank geht an Herrn Dr. G. Walter, der durch anregende Diskussionen und fördernde Kritik Wesentliches zum Entstehen dieser Arbeit beigetragen hat. Bei Herrn Dr. E. Scherzinger bedanke ich mich für die Einführung in die Methoden der Proteinexpression in *E. coli*, sowie für fördernde Anregungen und Diskussionen. Für verschiedene Anregungen zum Manuskript dieser Arbeit danke ich Dr. A. Lüking und Dr. C. Holz.

Contents

1. Introduction	1
1.1 Expression patterns	1
1.2 Protein expression	2
1.2.1 Expression of fusion proteins	3
1.2.2 Expression systems	4
1.3 Expression libraries to study protein interactions	6
1.4 Arrayed DNA libraries and high-density grids	6
1.5 Systematic protein expression	7
1.6 Protein characterisation by mass spectrometry	9
1.7 Antibodies	9
2. Objective	12
3. Materials	13
3.1 Laboratory equipment	13
3.2 Chemicals, nucleotides, antibodies and enzymes	14
3.3 Oligonucleotides	15
3.4 Kits	16
3.5 Other materials	16
3.6 Buffers and media	16
3.7 Strains	19
4. Methods	20
4.1 Plasmid constructs	20
4.2 PCR and DNA sequencing	20
4.3 Antibody affinity purification	21
4.4 cDNA library construction and arraying	22
4.4.1 Total RNA preparation	22
4.4.2 Selection of polyadenylated (poly(A) ⁺) RNA	22
4.4.3 cDNA synthesis	23
4.4.4 Gel electrophoresis of first strand cDNA	27
4.4.5 Vector digestion	27
4.4.6 Ligation	27
4.4.7 Preparing E. coli cells for electroporation	28

4.4.8 Transformation	28
4.4.9 Colony Picking	28
4.5 High-density filters for protein and DNA detection	29
4.5.1 DNA filters	29
4.5.2 Protein filters	29
4.6 DNA hybridisation screening of high-density filters	29
4.7 Antibody screening of high-density filters	31
4.8 Protein expression in <i>E. coli</i>	31
4.9 SDS-PAGE	31
4.10 Metal chelate affinity purification	32
4.10.1 Purification under denaturing conditions	32
4.10.2 Purification under native conditions	32
4.11 Tryptic digest	33
4.12 Protein expression and purification in microtitre plates	34
4.12.1 Protein expression	34
4.12.2 SDS-PAGE of whole cellular proteins	34
4.12.3 Metal chelate affinity purification	34
4.12.4 Solubility of expression products	35
4.13 Enzyme assays	35
4.13.1 GAPDH assay	35
4.13.2 Calmodulin assay	36
5. Results	38
5.1 Arrayed cDNA expression libraries	39
5.1.1 Expression vector pQE30NST	39
5.1.2 Construction of cDNA libraries	39
5.1.3 Preparation of high-density filters	42
5.1.4 Screening for recombinant protein expression	43
5.2 Rearranging of potential expression clones in the hEx1 library	46
5.2.1 Expression and purification in microtitre plates	46
5.2.2 Solubility	47
5.2.3 DNA sequence analysis	47
5.3 Identification of expression clones for specific genes	56
5.3.1 Screening of the hEx1 library with DNA probes	56
5.3.2 Biological activity of GAPDH and calmodulin	57
5.3.3 Detection of GAPDH and HSP90 α expression clones with antibodies and DNA probes	61
5.4 Characterisation of expression products by mass spectrometry	70
6. Discussion	74

6.1 Arrayed cDNA expression libraries	74
6.1.1 Robot technology and arrayed libraries	74
6.1.2 cDNA library construction	75
6.1.3 Detection of expression clones	77
6.2 Rearraying of the hEx1 library	78
6.2.1 Expression products of 96 random clones	78
6.2.2 DNA sequence analysis	78
6.2.3 Protein sizes predicted from DNA sequences and estimated from SDS-PAGE	79
6.3 Identification of specific expression clones	80
6.3.1 Screening for expression clones of nine human genes	80
6.3.2 GAPDH and HSP90 α expression clones	81
6.4 Characterisation of expression products	84
6.5 Perspectives	85
7. Abstract	88
7.1 Abstract in English	88
7.2 Zusammenfassung in deutscher Sprache	90
8. Appendix	93
8.1 Abbreviations	93
8.2 Curriculum vitae	96
8.3 Publications	96
9. References	97

7. Abstract

7.1 Abstract in English

Functional analysis of the proteins expressed in a specific tissue and/or stage of development is an essential step towards understanding biological processes. No technique has yet been available to go directly from sequence information on individual cDNA clones to the protein products. In the approach described here, arrayed cDNA expression libraries from human fetal brain (hEx1) and mouse adult kidney (mKd1) were established for the integration of DNA-based and protein-based experimental data. An *E. coli* plasmid vector (pQE30NST) for the expression of His₆-tag fusion proteins was used to clone cDNA synthesised by oligo(dT) priming. Using automated systems, 27,600 clones of the mKd1 and 193,500 clones of the hEx1 library were picked, grown in microtitre plates and arrayed on membrane filters at a density of 9,216 or 27,648 clones per filter. High-density DNA and protein filters were prepared and used to screen the cDNA libraries with DNA probes and antibodies in parallel. An antibody directed against the N-terminal tag sequence RGSH₆ of proteins expressed from pQE30NST (RGS-His antibody, Qiagen) was used to detect expression clones. Approximately 20% of clones in the mKd1 and hEx1 cDNA libraries were detected as putative expression clones by this antibody. 37,830 putative expression clones in the hEx1 library were combined and rearranged into new microtitre plates. A copy of this sub-library was given to the Resource Centre of the German Human Genome Project (RZPD, <http://www.rzpd.de>) to generate and distribute high-density DNA and protein filters. Expression products and DNA sequences of 96 randomly selected clones were analysed in detail. Protein expression and purification in microtitre plates was established. 68 out of 96 clones expressed proteins of at least 15 kd size (estimated from SDS-PAGE) which were detected by SDS-PAGE of whole cellular proteins or by metal affinity chromatography. DNA sequences derived from 5'-ends of cDNA inserts were obtained for 93 clones. 58 sequences matched human protein sequences in the SWISS-PROT and TrEMBL databases. 38 (66%) of these 58 sequences contained inserts cloned in the correct reading frame, i.e. the His₆-tag and a protein coding cDNA sequence were fused in frame. 66% of sequences matched to the beginning of a protein database sequence, and the

corresponding clones were therefore considered to contain complete protein-coding regions (full-length clones).

Expression clones for a panel of human genes were identified in the hEx1 library. Nine cDNA probes for BMP-7, calmodulin, COX4, GAPDH, hMSH2, HSP90 α , HSP90 β , RXR β and VDAC1 were used for screening by DNA hybridisation. Among the positive clones, putative expression clones were selected that were also detected by the RGS-His antibody on high-density filters. Expression clones for seven genes, but not for BMP-7 and VDAC1, were found with this strategy. GAPDH and calmodulin clones comprising full protein-coding regions expressing soluble His₆-tag fusion proteins were identified. Biological activity of His₆-tag GAPDH and calmodulin fusion proteins was shown with enzyme assays. DNA probes and the RGS-His antibody were used in combination to detect clones expressing GAPDH and HSP90 α fusion proteins. In parallel, specific antibodies directed against GAPDH and HSP90 α were used to identify expression clones on high-density protein filters. All clones identified by the protein-specific antibodies expressing His₆-tag GAPDH or HSP90 α fusion proteins were reliably detected by the RGS-His antibody. The RGS-His detected additional clones, which either contained GAPDH or HSP90 α inserts in an incorrect reading frame, or which expressed truncated proteins lacking the epitopes recognised by the protein-specific antibody. Estimated 90% of clones with inserts in an incorrect reading frame were not detected by the RGS-His antibody.

A protocol for the analysis of His₆-tag fusion proteins by matrix assisted laser ionisation/desorption mass spectrometry (MALDI-MS) was established. His₆-fusion proteins were purified under denaturing conditions by immobilisation on Ni²⁺-chelate-coated magnetic beads, followed by direct digestion with trypsin. The masses of the tryptic peptides were measured and compared to masses predicted from sequences in databases.

Arraying of cDNA expression libraries extends the application of DNA library arrays to protein-based screening techniques, thus creating a direct link between DNA-based and protein-based experimental data. The hEx1 expression library allows for the fast identification of expression clones for specific genes by DNA hybridisation or antibody screening. Proteins expressed from library clones can be directly used for functional assays, or, if inclusion bodies are formed, may be used as antigens to generate antibodies or possibly refolded *in vitro*. Libraries of expression clones displayed on high-density filters, in combination with high-

throughput protein expression and purification in microtitre plates, should be a feasible approach to the construction of gene product catalogues.

7.2 Zusammenfassung in deutscher Sprache

Die funktionelle Analyse der Proteine, die in bestimmten Geweben und/oder Entwicklungsstadien exprimiert werden, ist unverzichtbar für das Verständnis biologischer Vorgänge. Bis heute steht keine Technik zur Verfügung, um direkt von Sequenzinformation individueller cDNA Klone zu den Proteinprodukten zu gelangen. Bei der hier beschriebenen Vorgehensweise wurden geordnete cDNA Expressionsbanken aus humanem Fötalgewebe (hEx1) und Maus Nierengewebe (mKd1) hergestellt, um experimentelle Daten, die auf DNA-Basis oder Proteinbasis gewonnen wurden, integrieren zu können. Ein *E. coli* Plasmidvektor (pQE30NST) für die Expression His₆-markierter Fusionsproteine wurde für die Klonierung von cDNA benutzt, die mit einem oligo(dT) Primer synthetisiert wurde. 27.600 Klone der mKd1 und 193.000 Klone der hEx1 Bank wurden mit Hilfe von Robotern gepickt, in Mikrotiterplatten wachsen gelassen und auf Membranfiltern mit einer Dichte von 9.216 oder 27.648 Klonen pro Filtern angeordnet. Hochdichte-DNA und Proteinfilter wurden auf diese Art hergestellt und für das parallele Screening der cDNA-Banken mit DNA-Sonden und Antikörpern benutzt. Ein Antikörper gegen die N-terminale Sequenz RGSH₆ von Proteinen, die mit dem pQE30NST Vektor exprimiert wurden, wurde benutzt, um Expressionsklone zu identifizieren. Ungefähr 20% der Klone der mKd1 und hEx1 cDNA Banken wurden durch diesen Antikörper als mutmaßliche Expressionsklone identifiziert. 37.830 mutmaßliche Expressionsklone der hEx1 cDNA-Bank wurden zusammengefaßt und neu in Mikrotiterplatten angeordnet. Eine Kopie dieser neuen cDNA-Bank wurde dem Ressourcenzentrum im Deutschen Humangenomprojekt (RZPD, <http://www.rzpd.de>) für die Herstellung und Verteilung von Hoch-Dichte DNA- und Proteinfiltern gegeben. Die Expressionsprodukte und DNA-Sequenzen von 96 zufällig ausgewählten Klonen wurden analysiert. Zu diesem Zweck wurde die Expression und Reinigung von Proteinen in Mikrotiterplatten etabliert. 68 von 96 Klonen exprimierte Proteine von mindestens 15 kd Größe (bestimmt durch SDS-Gelelektrophorese), die nach SDS-Gelelektrophorese von Gesamtzellproteinen oder nach Metall-Affinitätschromatographie nachgewiesen wurden. cDNA-Inserts von 93 Klonen wurden vom 5'-Ende her ansequenziert. 58 dieser Sequenzen entsprachen Proteinsequenzen in der SWISS-PROT und TrEMBL Datenbank. 38 (66%)

dieser 58 Sequenzen enthielten Inserts, die im richtigen Leseraster kloniert waren, d.h. die His₆-Sequenz auf dem Vektor und eine Protein-kodierende cDNA Sequenz waren im richtigen Leseraster fusioniert. 66% war der Anteil der Sequenzen, die mit dem Anfang einer Proteinsequenz übereinstimmten, und von denen deshalb angenommen wurde, daß sie eine vollständige Protein-kodierende Sequenz enthalten („full-length Klone“).

In der hEx1 Bibliothek wurden Expressionsklone für eine Auswahl menschlicher Gene identifiziert. Neun cDNA-Sonden für BMP-7, Calmodulin, COX4, GAPDH, hMSH2, HSP90α, HSP90β, RXRβ und VDAC1 wurden für DNA-Hybridisierungs-Screenings benutzt. Unter den positiven Klonen wurden solche Klone, die auch durch den RGS-His-Antikörper detektiert wurden, als mutmaßliche Expressionsklone betrachtet. Expressionsklone für sieben Gene, jedoch nicht für BMP-7 und VDAC1, wurden so gefunden. Biologische Aktivität von His₆-GAPDH und Calmodulin Fusionsproteinen wurde durch Enzym-Assays nachgewiesen. DNA-Sonden und der RGS-His Antikörper wurden in Kombination eingesetzt, um Klone zu identifizieren, die GAPDH und HSP90α Fusionsproteine exprimieren. Parallel dazu wurden spezifische Antikörper gegen GAPDH und HSP90α eingesetzt, um Expressionsklone auf Hochdichte-Proteinfiltern zu identifizieren. Alle Klone, die von den Protein-spezifischen Antikörper erkannt wurden und His₆-GAPDH oder HSP90α Fusionsproteine exprimierten, wurde auch von dem RGS-His Antikörpern erkannt. Der RGS-His Antikörper detektierte weitere Klone, die GAPDH oder HSP90α Inserts enthielten, die entweder nicht im richtigen Leseraster kloniert waren, oder bei denen es sich um Teilesequenzen handelte, die nicht die von den Protein-spezifischen Antikörper erkannten Epitope umfaßten. Ungefähr 90% der Klone mit Inserts, die nicht im richtigen Leseraster kloniert waren, wurden durch den RGS-His Antikörper nicht erkannt.

Ein Protokoll für die Analyse von His₆-Fusionsproteinen durch *matrix assisted laser ionisation/desorption* Massenspektrometrie (MALDI-MS) wurde etabliert. His₆-Fusionsproteine wurden auf Ni²⁺-Chelator-beschichteten magnetisierbaren Partikeln immobilisiert und so unter denaturierenden Bedingungen gereinigt und direkt mit Trypsin verdaut. Die Massen der tryptischen Peptide wurden gemessen und mit vorhergesagten Massen aus Sequenzdatenbanken verglichen.

Die Verwendung von geordneten cDNA Expressionsbanken erweitert die Anwendung von geordneten DNA-Banken auf Protein-basierende Screening-Methoden, und schafft so eine direkte Verbindung zwischen experimentellen Daten, die auf DNA- und Proteinbasis

gewonnen werden. Die hEx1 Expressionbank ermöglicht eine schnelle Identifikation von Expressionsklonen für bestimmte Gene durch Screening mit DNA-Hybridisierungs-Sonden oder Antikörpern. Die Protein-Produkte von Klonen der cDNA-Bank können direkt für funktionelle Assays verwendet werden, oder können, falls *inclusion bodies* gebildet werden, als Antigene für die Produktion von Antikörpern verwendet oder möglicherweise *in vitro* gefaltet werden. Bibliotheken aus Expressionsklonen könnten, in Kombination mit Protein-Expression und -Reinigung in Mikrotiterplatten mit hohem Durchsatz, einen sinnvollen Weg zur Herstellung von Gen-Produkt-Katalogen darstellen.

8. Appendix

8.1 Abbreviations

2D-PAGE	two-dimensional polyacrylamide gel electrophoresis
A	adenine
A_{280}	absorbance at 280 nm
aa	amino acid
AMP	adenosine-5'-monophosphate
AP	alkaline phosphatase
APS	ammonium persulfate
ATP	adenosine-5'-triphosphate
BCIP	5-bromo-4-chloro-3-indolyl-phosphate
BMP	bone morphogenetic protein
bp	base pairs
BSA	bovine serum albumin
c	concentration
C	cytosine
CaM	calmodulin
CE	capillary electrophoresis
cAMP	3'-5'-cyclic AMP
cDNA	complementary DNA
Ci	Curie
COX4	cytochrome-c oxidase subunit IV
cpm	counts per minute
d	Dalton
dATP	deoxyadenosine-5'-triphosphate
DEPC	diethylpyrocarbonate
DIG	digoxigenin
DNA	deoxyribonucleic acid
DNase	deoxyribonuclease
dNTP	deoxyribonucleosid-5'-triphosphate
dpm	radioactive decays per minute
dsDNA	double stranded DNA
DTT	1,4-dithiothreitol
ϵ	absorption coefficient
<i>E. coli</i>	<i>Escherichia coli</i>
EDTA	ethylenediaminetetraacetic acid
EGTA	ethyleneglycol-bis-(2-aminoethyl ether)-N,N,N',N'-tetraacetic acid
ELISA	enzyme-linked immunoabsorbent assay
g	acceleration due to gravity, $g = 9.81 \text{ ms}^{-2}$
G	guanine
β -Gal	β -galactosidase, hydrolysis of lactose to glucose and galactose
GAPDH	glyceraldehyde-3-phosphate dehydrogenase
HEPES	N-(2-hydroxyethyl)piperazine-N'-2-(ethanesulfonic acid)
hEx1	human expression library 1
hIk-1	human ikaros 1 gene
His	L-Histidine
HSP	heat shock protein
IgG	immunoglobulin G
IPTG	isopropyl- β -D-thiogalactopyranosid
IMAGE	integrated molecular analysis of genomes and their expression
kbp	kilo base pairs
kd	kilo Dalton
<i>lacI</i>	<i>lac</i> repressor gene

<i>lacI</i> ^Q	<i>lacI</i> promoter mutation leading to enhanced expression
<i>lacZ</i>	β -galactosidase gene
LDH	lactate dehydrogenase
M	mole/litre
MALDI	matrix assisted laser desorption/ionisation
MALDI-MS	MALDI mass spectrometry
mKd1	mouse kidney library 1
MOPS	3-(N-morpholino)propanesulfonic acid
MPI-MG	Max-Planck-Institut für Molekulare Genetik, Berlin
mRNA	messenger RNA
MS	mass spectrometry
n.d.	not determined
NAD ⁺	β -nicotinamide adenine dinucleotide, oxidised
NADH	β -nicotinamide adenine dinucleotide, reduced
NBT	4-nitro blue tetrazolium chloride
NMR	nuclear magnetic resonance
NTA	nitrilotriacetic acid
OD ₆₀₀	optical density at 600 nm
oligo(dT)	oligo-deoxythymidine
PAGE	polyacrylamide gel electrophoresis
PBS	phosphate buffered saline
p	phosphate
PCR	polymerase chain reaction
PDE	phosphodiesterase
PE-C	pyridylethyl-cysteine
PEG	polyethyleneglycol
PK	pyruvate kinase
poly(A)	polyadenylic acid
PVDF	polyvinylidene difluoride
RBS	ribosome binding site
RGS	arginine-glycine-serine
RNA	ribonucleic acid
RNase	ribonuclease
rpm	revolutions per minute
RT	reverse transcriptase
RXR	retinoic acid X receptor
RZPD	Ressourcen-Zentrum Primär-Datenbank
SAP	shrimp alkaline phosphatase
SA-PMP	streptavidin paramagnetic particles
SDS	sodium dodecylsulfate
SDS-PAGE	SDS-polyacrylamide gel electrophoresis
SSC	saline sodium citrate
ssDNA	single stranded DNA
T	thymine
TAE	Tris-acetate-EDTA
Taq	DNA polymerase from <i>Thermus aquaticus</i>
TBE	Tris-borate-EDTA
TBS	Tris-buffered saline
TCA	trichloroacetic acid
TCEP-HCl	tris(2-carboxyethyl)-phosphine hydrochloride
TE	Tris-EDTA
TEA	triethanolamine
TEMED	N,N,N',N'-tetramethylmethylethylenediamine
TEN	Tris-EDTA-NaCl
TFA	trifluoroacetic acid
TOF	time of flight
Tris	tris(hydroxymethyl)aminomethan
tRNA	transfer RNA
U	units

URL	uniform resource locator
V	volume
VDAC1	voltage-dependent anion channel isoform 1
v/v	volume per volume
w/v	weight per volume
YAC	yeast artificial chromosome

8.2 Curriculum vitae

Name	Konrad Büssow
E-mail	buessow@mpimg-berlin-dahlem.mpg.de

Ausbildung

Sept. 1981 – Juni 1988	Schadow Schule Berlin, ABITUR
Okt. 1988 – März 1990	Technische Universität Berlin, Studiengang Biotechnologie
April 1990 – April 1995	Freie Universität Berlin, Studiengang Biochemie, DIPLOM BIOCHEMIE
13. Okt. 1994 – 13. April 1995	Diplomarbeit „Klonierung und Sequenzierung der Promotoren der DNA-Reparaturgene hMLH1 und hMSH2, Suche nach Polymorphismen“ am Wellcome Trust Centre for Human Genetics, Oxford
Mai 1995 – heute	Doktorand am Max-Planck-Institut für Molekulare Genetik Berlin, Abteilung Prof. Dr. H. Lehrach, Arbeitsgruppe Dr. G. Walter

8.3 Publications

Parts of this work have been published (129,130). A patent has been filed describing techniques established in this work (131).