

Genome-wide analysis of SPAK/OSR1 binding motifs

Eric Delpire and Kenneth B. E. Gagnon

Department of Anesthesiology, Vanderbilt University Medical Center, Nashville, Tennessee

Submitted 5 August 2006; accepted in final form 4 October 2006

Delpire E, Gagnon KB. Genome-wide analysis of SPAK/OSR1 binding motifs. *Physiol Genomics* 28: 223–231, 2007. First published October 10, 2006; doi:10.1152/physiolgenomics.00173.2006.—Based on the alignment of 12 sequences of protein motifs that interact with the kinases SPAK (Ste20-related proline alanine-rich kinase) and OSR1 (oxidative stress response 1), we performed genome-wide searches of the sequence [S/G/V]RFx[V/I]xx[V/I/T/S]xx, where x represents any amino acid. The “*Mus musculus*” search resulted in the identification of 131 mouse proteins containing 137 SPAK/OSR1 putative binding motifs. Similar numbers were found for human, zebrafish, fruit fly, and worm. A little more than half of the mouse proteins containing SPAK/OSR1 binding domains (53%) were also identified in the human search, whereas ~17–18% of these common hits were identified in the zebrafish search. The mouse proteins could be divided into two broad categories: 2/3 had an identified function, whereas 1/3 were either predicted or of unknown function. The known proteins were grouped as transport proteins, other membrane proteins, kinases, phosphatases, cytoskeletal, ribosomal, nuclear, enzymes, and others. Analysis of the location of the SPAK/OSR1 binding motif within the protein sequence revealed distribution throughout the entire length, but with preference to the extreme amino- or carboxyl termini for a large number of proteins. Analysis of the amino acid composition of the motifs revealed a preponderance of serine residues at positions 5, 6, 7, and 8. In summary, our new search found and thus confirms the 12 proteins previously shown to interact with the kinases and identifies 119 potential new targets for SPAK and OSR1 in the mouse proteome.

mouse genome; protein-protein interaction; docking site; Ste20 kinases; NCBI search; Ste20-related proline alanine-rich kinase; oxidative stress response-1

USING A LARGE-SCALE YEAST-2 hybrid screen and the amino terminus of KCC3 (K-Cl cotransporter-3), we uncovered a novel protein-protein interaction motif, recognized by two closely related Ste20 kinases: Ste20-related proline alanine-rich kinase (SPAK/PASK) and oxidative stress response-1 (OSR1) (5). Using deletion mutants and single residue mutagenesis, we determined that at minimum, nine residues in the target protein are required for kinase binding. The SPAK binding motif, published in 2002, has the sequence [R/K]F_x[V/I]xxxxx, where x represents any amino acid. Because the motif is based on yeast-2 hybrid data and this methodology cannot assess binding to residues located upstream of the positively charged arginine residue (that we will define as *position 1*), we cannot preclude the possibility that more than nine residues are involved in the interaction. Indeed, the smallest peptide (9 amino acids) that gives a positive interaction in a yeast-2 hybrid assay is fused to the binding domain of GAL4, and thus, amino acids linking this domain to

the peptide might in fact participate to the binding. Methodologies different from those using fusion proteins (different than yeast-2 hybrid and glutathione-S-transferase pull-down) are needed to resolve the minimum length required for SPAK/OSR1 binding. Subsequent studies from our laboratory and others have demonstrated binding of SPAK to several additional proteins (3–6). Sequence alignment of the motifs found in these proteins is presented in Fig. 1.

Several conclusions can be made from this alignment: First, serine or glycine residues are found preferentially at position –1, with the only exception being the original KCC3 target, which contains a valine residue at that position. Second, despite our previous demonstration of a conserved protein-protein interaction between SPAK and NKCC2 when we substituted a lysine for the arginine residue (2), of the 12 proteins listed in Fig. 1, only arginine residues are found at *position 1*. Third, amino acid substitution at *position 2* implies the necessity for a phenylalanine residue (2), a fact supported by the alignment represented in Fig. 1. Fourth, only valine or isoleucine residues are found at *position 4*. Finally, valine, isoleucine, threonine, and serine residues are found at *position 7*. Thus, based on this analysis, we propose the following refined motif: [S/G/V]RFx[V/I]xx[V/I/T/S]xx. Since proteins are made of 20 amino acids, at first approximation, the probability of finding one specific residue at any position is 1 divided by 20 or 0.05. Thus, the probability of finding the [S/G/V]RFx[V/I]xx[V/I/T/S]xx motif can be calculated to be 7.5×10^{-6} . Because some amino acids are more abundant than others, the actual probability of finding the motif increases slightly to 8.46×10^{-6} , indicating that if the SPAK binding motif was distributed within the genome by chance only, this motif should be found once every ~120,000 (118,151) residues.

To determine the number of SPAK/OSR1 binding motifs existing within a genome and identify the proteins containing this motif, we analyzed the National Center for Biotechnology Information (NCBI) protein database. Because the search had to accommodate multiple residues at any specific position, we created a small program written in Visual Basic that allowed us to identify the binding motifs within entire genomes. We provide here a complete list of mouse proteins containing [S/G/V]RFx[V/I]xx[V/I/T/S]xx motifs and provide information regarding conservation between different vertebrate genomes.

METHODS

The NCBI protein database (<http://www.ncbi.nlm.nih.gov/>) was searched for *Mus musculus*, and hits were saved in FASTA format within a single file. After opening the file and removing all “Hard Returns” using the find and replace function of WordPerfect, we saved the file in TEXT format. Next, using a small routine written in Visual Basic (Microsoft), we searched the entire text file for specific sequences allowing multiple residues per position. Basically, the first 60 characters are placed into 60 consecutive string variables, and the routine analyzes the file one variable at a time. Based on the identity

Article published online before print. See web site for date of publication (<http://physiolgenomics.physiology.org/>).

Address for reprint requests and other correspondence: E. Delpire, Dept. of Anesthesiology, Vanderbilt Univ. Medical Ctr., T-4202 Medical Center No., 1161 21st Ave. S., Nashville, TN 37232 (e-mail: eric.delpire@vanderbilt.edu).

	-1	1	2	3	4	5	6	7	8	9
KCC3a:	V	R	F	M	V	T	P	T	K	I
NKCC1:	S	R	F	Q	V	D	P	V	S	E
NKCC1:	G	R	F	R	V	N	F	V	D	P
NKCC2:	S	R	F	Q	V	H	V	I	N	E
WNK4 :	G	R	F	Q	V	T	S	S	K	E
WNK2 :	G	R	F	S	V	V	S	T	Q	D
AATYK:	S	R	F	T	V	S	P	T	P	A
AATYK:	S	R	F	S	I	T	H	I	S	D
HSP105:	G	R	F	V	V	Q	N	V	S	A
gelsolin:	G	R	F	V	I	E	E	V	P	G
CLH3 :	S	R	F	L	I	V	P	V	A	K
RELT :	G	R	F	R	V	A	R	I	P	E

Fig. 1. Alignment of the mouse KCC3 Ste20-related proline alanine-rich kinase/oxidative stress response-1 (SPAK/OSR1) binding motif with sequences of 11 other proteins experimentally shown to interact with the 2 kinases. Solid boxes delineate residues identified as essential for protein-protein interaction: R at position 1, F at position 2, and V/I at position 4. Dashed boxes identified residues exhibiting conservation at specific positions: V/S/G at position -1, and V/I/T/S at position 7.

of the first variable, a specific action is taken (Fig. 2), and the routine moves one character over, resetting all 60 variables until the end of the file is reached. All NCBI FASTA entries begin with a line of text describing the sequence (i.e., >gi 51592076 ref NP_032089.2 solute-carrierfamily37). Our program identifies new protein entries when the first variable = ">", the second = "g", and the third = "i". The routine then increments its protein counter by one and seeks the fourth "|" character, which signals the beginning of the name or description of the protein. The 60 characters following the fourth "|" character are saved temporarily in a string variable. Next, the routine seeks 30 consecutive capitalized letters, which signals the beginning of the sequence (this method was chosen since all NCBI entries do not start with a methionine), and starts incrementing its residue counter. When the routine identifies a serine, valine, or glycine at the first position

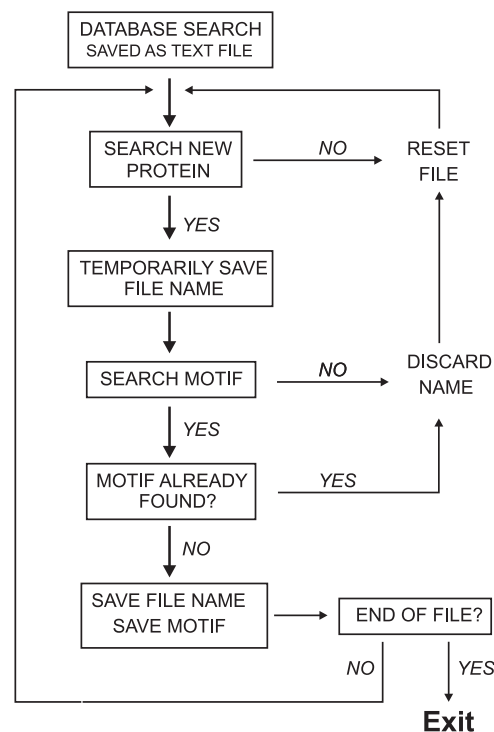


Fig. 2. Logic tree describing the general process used by the Visual Basic routine to identify SPAK/OSR1 binding motifs. See METHODS for details.

(first variable), followed by arginine and phenylalanine residues at positions 2 and 3, respectively, followed by a valine or isoleucine at position 5, and finally followed by valine, isoleucine, serine, or threonine residues at position 8, the routine captures the eight characters plus an additional seven characters following the motif into a temporary variable. This variable is then compared with 15 letter motifs previously saved, and if the string of 15 characters is unique, the motif belongs to a new protein and is saved along with the variable

Table 1. Summary of the SPAK/OSR1 binding motif search in the NCBI protein database

Search NCBI	Hits	Proteins with Motifs	Motifs Proteins	Duplicate	Residues	Time
All mouse	189,148	170	175	781	75,235,204	4h
Receptor	21,389	32	33	67	8,474,634	28'
Kinase	8,228	25	30	36	4,819,562	21'
Phosphatase	2,412	2	2	14	1,238,240	4'
Enzyme	3,018	11	13	12	1,509,775	5'
Cytoskeleton	1,389	9	9	0	1,079,907	3'
Transport	3,331	15	18	22	1,970,996	6'
Channels	2,601	9	9	41	1,551,394	5'
K ⁺ channels	551	1	1	3	279,708	57"
Na ⁺ channels	272	3	3	18	212,985	37"
Cl ⁻ channels	184	0	0	0	113,045	20"
Ca ²⁺ channels	307	0	0	0	267,320	45"
Pumps	147	1	1	0	82,138	16"
<i>Homo sapiens</i>	330,694	215	219	845	115,230,700	6 h 46'
<i>Danio rerio</i>	51,687	136	138	156	23,297,462	1 h 19'
<i>Drosophila</i>	81,930	103	104	275	39,321,787	1 h 56'
<i>C. elegans</i>	55,750	167	168	145	27,179,094	1 h 22'

The top portion of column 1 lists the results of several Ste20-related proline alanine-rich kinase/oxidative stress response-1 (SPAK/OSR1) motif searches using different criteria attached to "*Mus musculus*." The bottom portion of column 1 lists searches made for different species. The 2nd column lists the number of NCBI hits attached to each search. The 3rd column reports the number of proteins identified by the VisualBasic routine containing SPAK/OSR1 binding motifs. The 4th column reports the total number of motifs, as some proteins contain more than one motif. The 5th column lists the number of proteins that were eliminated because they were repeat hits. The 6th column reports the total number of amino acid residues analyzed by the routine, and the last column reports the times taken by the VisualBasic routine to return the data.

Table 2. Mouse proteins with SPAK binding domains

Protein	Motif	Accession	Position, %	Conserv.
<i>Transport Proteins</i>				
Fatty acid (Slc27a4)	SRFHIPQVAE	NP_036119	58	H
Na-K-2C1 cotransporter 1, NKCC1 (Slc12a2)	SRFQVDPVSE	NP_033220	6.3	H
	GRFRVNFVD		10.1	
Na-K-2C1 cotransporter 2, NKCC2 (Slc12a1)	SRFQVHVINE	P55014	1.4	H
K-C1 cotransporter 3, KCC3 (Slc12a6)	VRFMVTPTKI	NP_598410	1.1	H
Organic cation transporter	GRFQVRLTIL	AAH21449	2.7	
ATP-binding cassette subfamily C member 9 isoform b	VRFAVKAIIS	NP_066378	38.1	H
Nicotinamide nucleotide transhydrogenase	VRFGIHPVAG	NP_032736	88.2	H
Gap junction membrane channel protein alpha-1	VRFWVLQIIF	NP_034418	19.6	
Gap junction membrane channel protein alpha-7	VRFWVFQIIL	NP_032148	18.7	HZ
Gap junction membrane channel protein alpha-9	SRFYIIQVVF	NP_034420	61.1	HZ
Gap junction alpha-12 protein (connexin-47)	VRFWVFQIVV	Q8BQU6	17.5	HZ
<i>Other Membrane Proteins</i>				
SH3/ankyrin domain gene 3	VRFVVRVSVA	NP_067398	93.8	H
Transmembrane protein 24	VRFPISVTQQ	NP_082185	22.8	H
Polyductin (PKHD1 gene)	VRFYVEGSNI	NP_694819	32.2	
Taste receptor type 2 member 136	VRFCICVWVI	NP_851793	30.1	
Olfactory receptor MOR218-3	GRFKIFSTCI	AAL61334	74.6	
Down syndrome cell adhesion molecule	SRFLITSTGA	NP_112451	8.5	H
Interleukin 12 receptor, beta 2	SRFIVRVTAI	NP_032380	24	
Tumor necrosis factor receptor superfamily, member 19-like	GRFRVARIPE	AAI00311	81.2	HZ
<i>Kinases</i>				
Aortic preferentially expressed gene 1 (S/T kinase) <i>alternatively spliced exon</i>	VRFRVACSNR	NP_031489	72.9	
	GRFGVVRSC	BAC65770	78.9	
Serine/threonine kinase SADB	VRFQVDISS	NP_001003920	61.3	H
Brain apoptosis-associated tyrosine kinase, AATYK	SRFTVSPTPA	NP_031403	97	HZ
	SRFSITHIS		98	
WNK1	GRFQVSVTMD	NP_941992	78	H
	GRFQVTTTA		82	
	GRFSVSRT		82	
WNK2	GRFSVVSTQD	AAH46464	73.7	H
WNK3	GRFQVITVPQ	XP_986322	74.7	HZ
WNK4	GRFQVTSSKE	NP_783569	81.4	H
<i>Phosphatases</i>				
Protein phosphatase 2 regulatory subunit B (59 kDa)	VRFTVNPSDP	XP_923875	2.6	
Protein phosphatase 2 regulatory subunit A	VRFNVAKSLQ	NP_058587	91	HZ
<i>Cytoskeleton</i>				
Gelsolin	GRFVIEEVP	NP_666232	86.5	HZ
Titin isoform N2-B	VRFGVTITVH	NP_082280	93.6	HZ
	GRFHIENTD		99.8	
Dynein axonemal heavy chain 11	VRFSINDSTT	NP_034190	31.2	
Nexilin (actin filament-binding protein)	VRFTVKVTGE	NP_955759	87.3	HZ
Scin protein	GRFIIIEVPG	AAH63328	87.3	
Advillin	GRFLVTEVTD	NP_033765	76.4	
<i>Ribosomal</i>				
Ribosomal protein S11	VRFNVLKVTK	NP_038753	87.3	HZ
Ribosomal protein S23	VRFKVVKVAN	AAI00604	82.5	HZ
Mrps25 protein (mitochondrial)	GRFPIRRTLQ	AAH22953	2.9	HZ
<i>Nuclear</i>				
Structure-specific recognition protein	VRFYVPPTQE	NP_892035	22.9	HZ
RecQ protein-like (helicase)	VRFVIHHSMS	NP_075529	59.1	HZ
Bloom syndrome protein homolog (helicase)	VRFVIHASLP	NP_031576	68.2	H
DNA-dependent ATPase SNF2L	VRFEVSPSYV	AAK52453	16.4	
Cardiac-specific RNA helicase	SRFRIITTC	AAK77049	40	H
Nesprin (synaptic nuclear envelope 1 isoform 2)	SRFQIQQTAN	NP_071310	51.6	H
RanBP2 protein	SRFKVSRIGK	CAA60778	0.3	
Bromodomain containing 1	SRFGISTSLE	NP_001028446	98.8	
Minichromosome maintenance-deficient 2 mitotin	SRFDVLCVVR	NP_032590	72.5	H
Aquarius (helicase)	SRFAIRKIML	NP_033832	3.8	H
TAR DNA binding protein isoform 3	GRFGVHLISN	NP_001003899	92.9	H

Continued

Table 2.—Continued

Protein	Motif	Accession	Position, %	Conserv.
<i>Nuclear (Continued)</i>				
Myst4 protein	GRFLIDFSYL	AAH95974	80	HZ
Zmym4 protein	SRFMIELTKL	AAH50924	82.2	H
Pseudouridine synthase (Rpusd3)	SRFHVMATGR	AAH3808	60.7	
<i>Enzymes</i>				
Glycosylphosphatidylinositol specific phospholipase D1	GRFQIGRVYI	NP_032182	48.5	
Similar to O-acyltransferase (membrane bound)	VRFYITLSSL	XP_134120	36.7	
Acetyl-coenzyme A carboxylase alpha	SRFIIGSVSE	NP_579938	1	H
Endothelin converting enzyme 1	SRFRVIGSLS	NP_955011	95.3	HZ
Phosphodiesterase 1b (Pde1b)	SRFKIPTVFL	AAH58531	42.4	H
Phosphodiesterase 1c isoform a	SRFKIPISAL	NP_035184	30.6	H
Neurolysin (metallopeptidase M3 family)	SRFDIEMSMR	NP_083723	18.5	H
<i>Extracellular Space</i>				
Cathelicidin (antimicrobial peptide)	VRFRVKETVC	NP_034051	42.8	
Roundabout homolog 1	SRFSVSQTGD	NP_062286	22.6	HZ
Pentraxin (Crp, membrane-bound calcium-binding protein)	VRFMVSEIPE	NP_031794	44.9	
Fibroblast growth factor 4	SRFFVAMSSR	NP_034332	63.9	H
<i>Other Proteins</i>				
Dennd2c protein (suppression of tumorigenicity)	VRFFVELVGH	AAH46440	79.1	H
Dennd2a gene	VRFFVEIVGH	NP_766065	89.7	
Selenoprotein X 1	SRFUIFSSSL	NP_038787	79.3	HZ
Heat shock protein 105	GRFVVQNVSA	NP_038587	53.6	HZ
Eef2 protein (elongation factor)	VRFDVHDVTL	AAH60707	81.3	H
High glucose-regulated protein 8	GRFDVRWIFV	NP_663368	85.1	HZ
Procollagen type XVIII alpha 1	GRFGINGSYA	NP_034059	36.7	
IgH antibody heavy chain VDJ region	GRFIVRSRD	AAB95229	54.8	
Follistatin-like 4	GRFIVSVSNK	NP_796033	82.9	
Flywch1 protein	GRFLVYESFL	AAH25645	60.8	H
FERM domain containing 4A	VRFYIESISY	NP_766063	8.1	HZ
Golgi-associated band 4.1-like protein	VRFYIESISF	NP_660130	8.6	H
Carcinoembryonic antigen	VRFHVHPILL	CAA33409	29.7	
Low-density lipoprotein receptor-related protein 6	SRFVIINTEI	NP_032540	10.7	
Proline-rich 14 protein	SRFIRRTPV	BC006909	58.5	H
Poly(ADP-ribose) polymerase family, member 14	SRFPVDVVVN	NP_001034619	45.3	
Ccdc60 protein	SRFLIQCVKI	AAH75662	27.5	H
HEF-like protein	SRFLIPRVEQ	NP_001028710	47.4	
Sfrs14 protein	SRFGIEIKW	AAH57305	44.9	
Ig heavy chain var. region MANEX1B against dystrophin	GRFTISRVRN	S56008	60.9	H
Ig heavy chain variable region scfv fragment	GRFTISKTSS	AAL15163	53.9	
IgH antibody heavy chain VDJ region	GRFFVRSRD	AAB94782	56.7	
Polymeric immunoglobulin receptor	GRFSVLITGL	NP_035212	39.6	
Pdxx protein	GRFSVELTRG	AAH99459	38.6	
Gtf3c1 protein	VRFRISNSST	AAH32208	5.7	H
Glutaredoxin cysteine-rich 1 protein	VRFRIASSHS	NP_001018019	5.5	H
D19Ert652e protein	VRFQVNLTVQ	AA107403	61.3	
Vacuolar protein sorting 37D	GRFGILSTGQ	NP_808242	9.7	
Msid2	GRFKIDVSDT	AAS87379	59.8	
Left right determination factor 1	GRFLVSETST	NP_034224	23.8	
<i>Unknown Proteins</i>				
Unnamed protein product	VRFLICFIGM	BAC33608	94.5	
Unnamed protein product	VRFLVSWSWL	BAE26787	95.3	
Unnamed protein product	VRFCVFPTGG	BAE24015	56.3	
Unnamed protein product	SRFTVAAIRR	BAB25824	74.2	
Unnamed protein product	SRFSIISSSS	BAE38028	3.9	H
Unnamed protein product	SRFRVKTLE	BAC27063	22.9	H
Unnamed protein product	SRFLVILSFF	BAE33644	32.4	
Unnamed protein product	GRFSVCFSFP	BAE24254	90	
Unnamed protein product	GRFFIPSVTQ	BAE23369	41.9	
Unnamed protein product	GRFKVFRISP	BAE42406	78.9	
Unknown protein	VRFSVPGTTW	BAE24216	36.1	
Unknown protein	VRFRIAKIVA	BAE22912	38.1	
RIKEN cDNA 1110057K04	SRFPVWIISH	AAH46986	22.4	
RIKEN cDNA 4631403P03	GRFRVTHIEK	AK044097	55.1	H

Continued

Table 2.—Continued

Protein	Motif	Accession	Position, %	Conserv.
<i>Unknown Proteins (Continued)</i>				
C730048C13Rik protein	GRF QILQISF	AAH50265	3.1	
mKIAA0539 protein	GR F TVNKVAL	BAD90144	72.4	
Amphoterin-induced gene	SRFPV E LKID	AAH10598	34.8	
<i>Predicted Proteins</i>				
Similar to 40S ribosomal protein S4, X isoform	GRF AVHHITL	XP_001004419	40.1	H
Similar to TBC1 domain family member 4 (Akt substrate)	GRF EINLISP	XP_619279	45.5	
Sodium channel voltage-gated type IX alpha 8	SRFSVSQVAN	XP_904778	68.1	
Importin 9 isoform 5 (1, 6)	SRF TVAMSPE	XP_903902	72.4	H
Similar to hyaluronoglucosaminidase 4	SRF RVRESLR	XP_995969	57.8	
Similar to proline-rich 6 isoform 4 (or 6)	VRF EVWASAD	XP_921840	53.2	H
Similar to Bex4 protein	GRF VVQGTEV	XP_982566	67	
Similar to T-cell receptor alpha chain V region HPB-MLT precursor isoform 2	GRF SVKHSKA	XP_925983	51.3	H
Hypothetical protein LOC332309	VRF KISSSYS	NP_001028598	7.5	H
Hypothetical protein LOC382137	VRF GVDTQL	NP_941077	9.8	H
Hypothetical protein LOC107373 (hydrolase, peptidase)	VRF CIHAVGS	NP_080916	29.2	
Hypothetical protein LOC231279	SRF PVEDIRN	NP_766299	6.6	H
Hypothetical protein LOC71874	SRF IVDYSRD	NP_765999	13.3	
Hypothetical protein LOC212943 isoform 1	SRF FIDFSDI	XP_135029	71.5	H
Hypothetical protein LOC73453	GRF SISPSHD	NP_898920	42.3	
Hypothetical protein LOC100532	GRF RVTKVEH	NP_666035	77.7	H
Hypothetical protein LOC75619	GRF SILNSQH	NP_766010	47	
Hypothetical protein LOC226089 isoform 3	GRF LIIFIAI	XP_905557	92.8	
Hypothetical protein LOC75905	GRF LVSTSTR	NP_780437	56	
Hypothetical protein LOC210463	GRF QICMIL	NP_759010	2.5	
Hypothetical protein LOC271981 (kinase)	GRF QILKTIT	NP_766620	5.2	H
Similar to Kelch repeat and BTB domain-containing protein	SRFQIPSVFT	XP_915191	18.3	H

Column 1 lists all mouse proteins that were identified as including [S/G/V]RFx[V/I]xx[V/I/T/S]xx binding sites. The exact sequence of the motif(s) is reported in the 2nd column with residues at positions 1, 2, and 4 highlighted (bold). The 3rd column contains the accession number of the protein, the 4th column indicates the location of the motif within the protein (as a percentage). In the last column, we indicate whether the motif is also found in the human (H) or zebrafish (Z) orthologs.

containing the name of the protein, and the motif counter is incremented by one. In contrast, if the string of 15 character is already saved, this indicates a redundant protein and the “duplicate protein” counter is incremented by one. The likelihood that two different proteins contain the exact same 15 characters is infinitesimal, as the

probability to find a specific string of 15 residues can be estimated at 3×10^{-20} . The routine then continues its search without saving the motif. At the end of the search, all variables (protein counter, motif counter, protein names, and motifs) are copied into a new single text file.

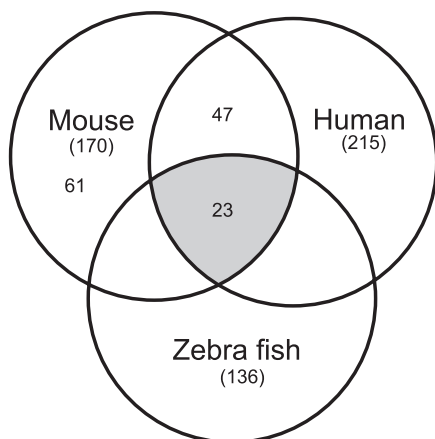


Fig. 3. Conservation of SPAK/OSR1 binding motifs. The value in parentheses indicates the number of protein hits in 3 vertebrate genomes: mouse (*Mus musculus*), human (*Homo sapiens*), and zebrafish (*Danio rerio*). We identified a total of 131 mouse proteins from the 170 hits. There were 23 out of 131 proteins conserved in mouse, human, and zebrafish genomes. Sixty-one proteins were specific to the mouse genome, and 70 (47 + 23) proteins were conserved between the mouse and human genomes. The list of these proteins is reported in Table 2.

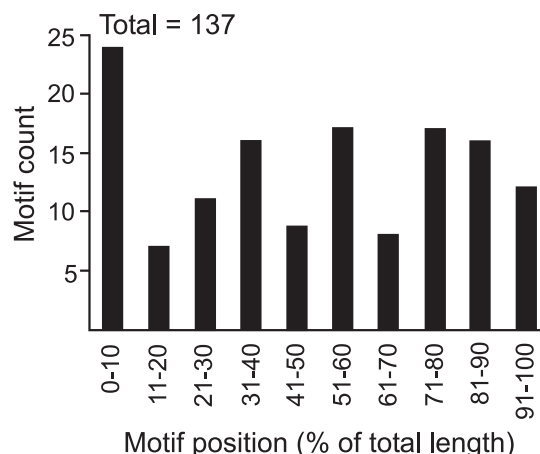


Fig. 4. Distribution of the SPAK/OSR1 motifs within the length of proteins. For each motif identified, its location (*position -1*) was divided by the total length of the protein. There were 137 motifs identified in the mouse search. Note that the motif is distributed throughout the length of the protein with a larger number of proteins with SPAK/OSR1 binding motif located within the 1st 10% of the protein length. The motif location for each individual protein is reported in Table 2.

Table 3. Mouse proteins with putative [S/G/V]KFX[VI]xx[VI/T/S]xx motifs

Protein	Motif	Accession	Position (%)
<i>Transport Proteins</i>			
Aquaporin 6	GKFAVHWIFW	NP_780296	70
Plasma membrane calcium ATPase 1	VKFFIHGVTV	NP_080758	34.2
Potassium inwardly rectifying channel J15 isoform a	VKFHVDSSSE	NP_001034146	65.2
Equilibrative nucleoside transporter 3	SKFFVPLTVF	Q99P65	72.7
ATP-binding cassette, subfamily A (ABC1), member 17	SKFGIGYSLQ	NP_001026792	91.6
Sperm-associated cation channel 2	SKFIIISLIFL	NP_694715	18.5
Solute carrier family 12, member 2	SKFRIDFSDI	NP_033220	88.8
<i>Other Membrane Proteins</i>			
Purinergic receptor P2X1	GKFDIIPMT	NP_032797	81.2
Purinergic receptor P2X3	GKFNIIPTII	NP_663501	79.1
Purinergic receptor P2X4	GKFDIIPMTI	NP_035156	83.4
Purinergic receptor P2X5	GKFSIIPTVI	NP_201578	72.3
Cadherin 23	GKFEIDESTG	NM_023370	37.4
Olfactory receptor 1496	VKFFVGLSVI	NP_667200	63.1
Olfactory receptor 898	VKFAVKRTMK	NP_667082	95.9
Olfactory receptor 116	SKFTVRIICP	NP_666843	96
Olfactory receptor 452	GKFSVFTSKL	NP_001011869	96.5
Huntingtin	VKFFVMTVEA	NP_034544	71.5
Protocadherin alpha 5	SKFTIDSSSG	NP_034089	31.1
Cadherin-related neural receptor 6	SKFTIDSSSG	BAA29051	28.4
Killer activatory receptor-like protein p91D	GKFFIPSVTQ	AAC40073	41.9
PTX1 protein isoform 1	GKFFIVEIICC	NP_080444	89.1
Macrophage scavenger receptor 2	SKFTISAISK	NP_109632	31.8
Anti-ricin antibody light chain variable region	VKFLIYSTSN	ABD47464	40
<i>Kinases</i>			
Protein kinase cGMP-dependent type II	SKFCVACVTE	AA113206	72.6
Brain-selective kinase 2 isoform (alpha, beta, gamma isoforms)	VKQVDITYT	NP_083702	90
Bruton agammaglobulinemia tyrosine kinase	SKFPVRWSPP	NP_038510	84.5
Thymidylate kinase family LPS-inducible member	GKQVIAIEG	NP_065582	50.4
Hexokinase domain containing 1	SKFRVLKVQV	NP_663394	9.6
<i>Phosphatases</i>			
Protein tyrosine phosphatase nonreceptor type 9	GKFTILNVRD	NP_062625	15.9
<i>Nuclear</i>			
Topoisomerase (DNA) II alpha	VKFKVIMTEE	NP_035753	63.2
Topoisomerase (DNA) II beta	VKFFVVKMTEE	NP_033435	60.5
DNA methyltransferase 3A isoforms 1 & 2	GKFSVVCVEK	NP_031898 (1)	36.3
Polymerase (RNA) II (DNA directed) polypeptide C	VKFIIENTDL	NP_033116	6.9
Methyl-CpG binding domain protein 4	GKFDVYFISP	NP_034904	16.4
MYST histone acetyltransferase 1	GKFLIAFSYE	NP_080646	71.8
SWI/SNF-related matrix associated actin dependent	GKFNVLTTY	NP_035547	52.7
Regulator of chromatin, subfamily a, member 4			
Chromodomain helicase DNA binding protein 3	VKFKVLLTSY	CAI35246	44
DEAH (Asp-Glu-Ala-His) box polypeptide 35	VKQVQPKVSS	NP_665685	86.9
<i>Enzymes</i>			
RAB5A-B-C, members RAS oncogene family	VKFEIWDTAG	NP_080163 (5A)	32.1
Dipeptidylpeptidase 4	VKFFIVNIDS	NP_034204	34.2
2'-5' Oligoadenylate synthetase 1G	VKFEVQSSWW	NP_035982	32.4
2'-5' Oligoadenylate synthetase 1C & 1F	VKFEVQSSEE	NP_291019 (1C)	32.9
Carnosine dipeptidase 1 (metallopeptidase M20 family)	GKFSIRLVPT	NP_803233	70.3
Carnosine dipeptidase 2 (metallopeptidase M20 family)	GKFSIRLVPD	NP_075638	71.2
Arhgef10 (guanine nucleotide exchange factor for Rho/Rac/Cdc42-like GTPases)	VKFIVSATAF	AAH66074	85.8
RAS p21 protein activator 2	SKFTIEDSVA	NP_444498	93.7
ADP-ribosylation factor 6	VKFNVDVGG	NP_031507	32.6
Phosphoglucosyltransferase 1	GKFEISAIRD	NP_079976	83.7
Transglutaminase 3 E polypeptide	GKFKVTGILA	NP_033400	70
4-Hydroxyphenylpyruvic acid dioxygenase	VKFAVLQTYG	NP_032303	33.3
Peptidase 4	SKFNVNNTIL	NP_032846	33.9
Rab geranylgeranyl transferase a subunit	SKFLVENSVL	NP_062392	75.3
UTP-glucose-1-phosphate uridylyltransferase	SKFKIFNTNN	AAH25585	62.6
UDP-glucuronosyltransferase 2 family, member 5	SKFDVLLSDP	BAE22997	28.7
UDP-glucuronosyltransferase 2 family, polypeptide B36	SKFDVILSDA	NP_001025038	27.2
Serine (or cysteine) peptidase inhibitor, clade A, member 3B	GKFIIVDRSRH	NP_766612	54

Continued

Table 3.—Continued

Protein	Motif	Accession	Position (%)
<i>Extracellular Space</i>			
Tubulin tyrosine ligase-like family, member 12	VKFDIRYIVL	NP_898838	72.3
Pregnancy zone protein	VKFRVVSVDI	NP_031402	9.5
<i>Other Proteins</i>			
Pad-I-like isoform 1	VKFTISTTSK	NP_081569	68
YTH domain family 1	GKFDVKWIFV	NP_776122	84.4
Ythdf3 protein	GKFEVKWIFV	AAH67040	85.3
Sulfatase modifying factor 1	GKFPVSNTGE	NP_666049	60.9
Ly-6C variant	GKFPVYIYIK	BAA13049	61.8
Serine (or cysteine) proteinase inhibitor, clade B (ovalbumin)	GKFRIGFIDE	NP_082247	58.2
Neuron specific gene family member 2	GKFRVPKIAE	NP_032767	35.1
Phf13 protein	SKFDIRRSNR	AAH59282	93.2
Nup214 protein	VKFAVQDVND	AAH39282	16.5
Sperm acrosome associated 1	VKFTVYTTNE	NP_080569	67.2
Translocon-associated protein alpha	VKFLVGFTNK	AAK82421	35.3
Tumor-associated antigen	VKFLVAKVYE	AAB17502	23.9
Hydin (hydrocephalus-inducing protein)	GKFTISPSI	AAO44953	61.7
Eukaryotic elongation factor, selenocysteine-tRNA-specific	SKFKIHITIQ	NP_055547	87.3
DEP domain containing 1B	SKFIIHNVYS	NP_848798	40.2
Apolipoprotein B	SKFKVSHVEK	AAH38263	97
Immunoglobulin H-chain V-region	SKFDVWGTGT	BAD24599	87.5
Melanoma-derived leucine zipper, extra-nuclear factor	SKFAIKGIIN	NP_113555	46.5
<i>Unknown Proteins</i>			
Unnamed protein product	VKFCIDASQP	BAC28336	2.8
<i>Predicted Proteins</i>			
Hypothetical protein LOC67105	VKFYIEGSEP	NP_080245	78.3
Similar to Serum amyloid P-component precursor	GKFDVKQSFV	XP_136329	70.6
Similar to Extracellular calcium-sensing receptor	GKFLVGIIGA	XP_981998	17.3
Similar to Zinc finger protein 208	SKFFVYPSRL	XP_995677	66.8
Similar to FLJ44048 protein	VKFTVPVTLT	XP_918762	22.2
Similar to 60S ribosomal protein L38	VKFRVRCSTRY	XP_983985	45.7
Similar to trinucleotide repeat containing 6A	SKFVVGSSSN	XP_987987	17.9
Similar to enolase 1, alpha	SKFGVNILG	XP_357379	22.1
Similar to mucin 17	SKFGIATSKN	XP_992907	71
Similar to Zinc finger protein 208	SKFFVPHSQL	XP_995677	66.8
Similar to nucleolar protein NOP56	VKFNVNRVDN	BAB27647	45.4
Similar to CG31901-PA	GKFFVISTSDG	XP_913556	80.4
Similar to GTPase activating protein testicular GAPI	VKFLIENSLK	XP_995294	44.7
Similar to putative binding protein 7a5	GKFTVSDIRK	XP_995079	47.7
Hypothetical protein LOC210711	GKFTVNLGSGC	NP_666067	60.3
Hypothetical protein LOC241134	GKFTVPASHS	NP_766437	23.7
Hypothetical protein LOC329562	GKIFYVAIVCS	NP_808519	76.2
Hypothetical protein LOC67484	VKFEICVSSK	NP_080465	16.5
Hypothetical protein LOC66367	VKFAIKRTLK	NP_780316	90.1
Hypothetical protein LOC268706	SKFNLGTVS	NP_848861	55.5
Hypothetical protein LOC626359	VKFEVMCVVL	NP_001033016	70.8
Hypothetical protein from cochlea	VKFQVVKIKG	CAA67982	92.2

Column 1 lists all mouse proteins that were identified substituting the arginine residue at *position 1* with a lysine residue. The sequence of the alternative SPAK/OSR1 motif is reported in the 2nd column. The 3rd column lists the accession number, and the 4th column indicates the location of the motif within the protein (as a percentage).

For each of the proteins returned by the Visual Basic routine, we queried the NCBI protein database to locate and confirm the motif within the amino acid sequence of the protein. Additionally, the curation process provided the opportunity to eliminate any remaining redundant or nonmurine proteins and allowed grouping of the proteins into categories based on their general description and features.

RESULTS AND DISCUSSION

The NCBI protein database (<http://www.ncbi.nlm.nih.gov/>) was first searched for *Mus musculus*. The search was performed on June 14, 2006, and returned some 189,153 hits.

More directed searches such as *Mus musculus* AND “descriptor” returned: 21,389 hits for “receptors,” 8,228 hits for “kinases,” 3,018 hits for “enzymes,” 1,389 hits for “cytoskeleton,” 3,331 hits for “transport,” 2,412 hits for “phosphatases,” 2,601 hits for “channels.” Further defining the type of ion channel reduced the number of hits to 551 for “K⁺ channels,” 307 for “Ca²⁺ channels,” 272 hits for “Na⁺ channels,” and 184 hits for “Cl⁻ channels” (see Table 1). An obvious first observation that can be made from these searches is that NCBI returns many more hits than the number of proteins existing in any one specific genome. Indeed, the NCBI search of *Mus musculus*

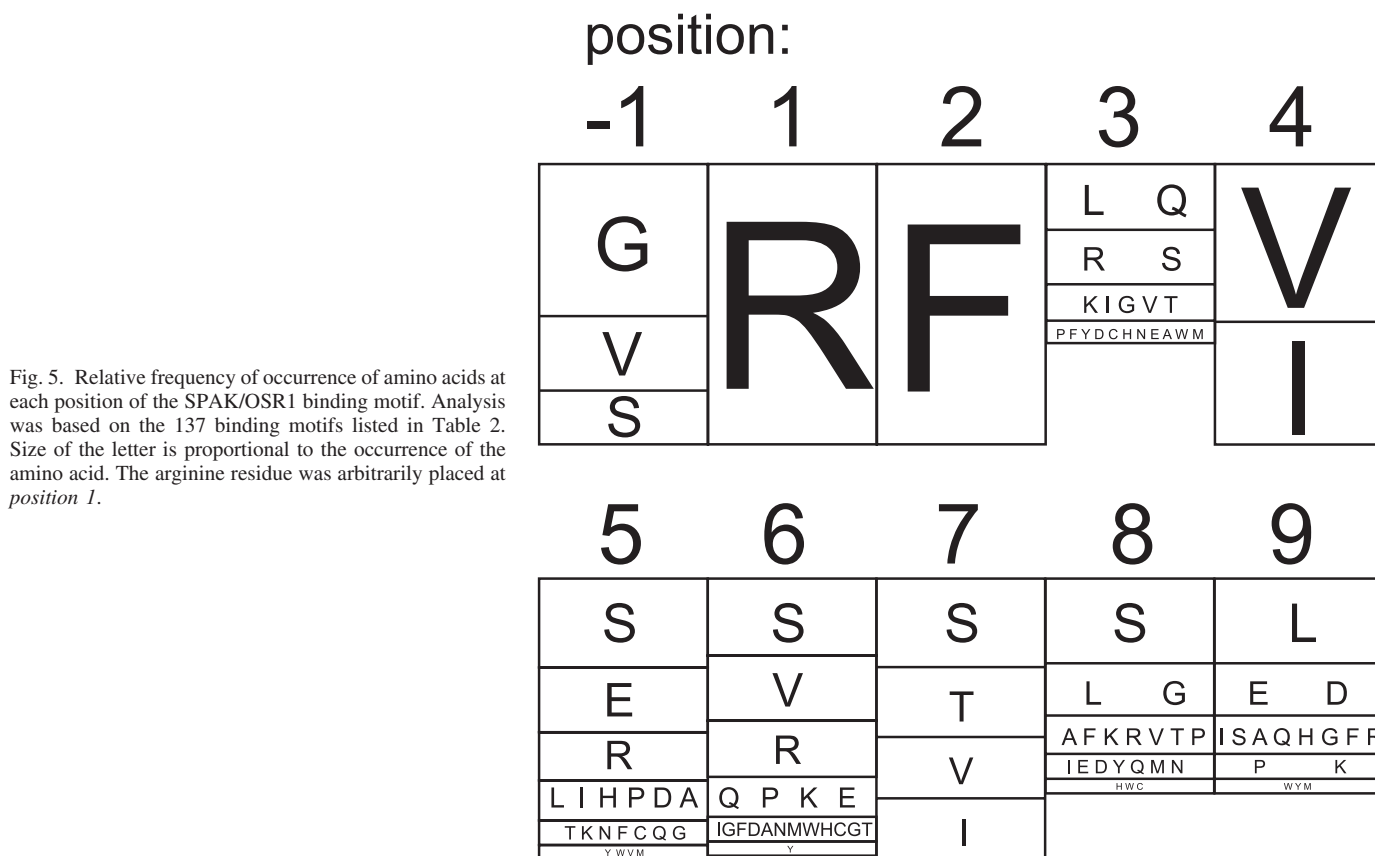


Fig. 5. Relative frequency of occurrence of amino acids at each position of the SPAK/OSR1 binding motif. Analysis was based on the 137 binding motifs listed in Table 2. Size of the letter is proportional to the occurrence of the amino acid. The arginine residue was arbitrarily placed at position 1.

returned close to 200,000 hits, whereas estimation of the number of proteins in the mouse genome is closer to 24,000–25,000 (7). A second observation is that, due to the extensive cross-referencing of NCBI entries, proteins are often identified using unrelated or indirect searches. As an example, only a fraction of the 8,228 entries found with *Mus musculus* AND “kinases” are actual kinases. The Na-K-2Cl cotransporter (NKCC1) is identified in the *Mus musculus* AND “kinases” search due to the multiple NCBI links existing between SPAK/PASK kinases and the cotransporter.

Using a routine written in Visual Basic, we searched for occurrence of the [S/G/V]RFx[V/I]xx[V/I/T/S]xx within the 189,153 downloaded mouse proteins. The routine identified 170 proteins containing 175 motifs. Analysis of the proteins showed that the routine was able to eliminate all redundant entries and alternatively spliced products, thus identifying one product per gene. However, a small number of proteins identified in the “*Mus musculus*” search were not mouse proteins but belonged to other species. After curation and elimination of these proteins the final list comprised 131 proteins containing 137 motifs. There were five instances where a single open reading frame encoded a protein containing more than one SPAK/OSR1 binding motifs. WNK1 contains three independent SPAK/OSR1 binding motifs, and four other proteins (NKCC1, AATYK, titin isoform N2B, and mKIAA1297) contain two SPAK/OSR1 motifs. The complete list is presented in Table 2. Two-thirds of these proteins are involved in transport, enzymatic, cytoskeletal, nuclear, ribosomal, and other functions, whereas the remaining one-third of the proteins listed in

Table 2 either are only predicted or have undetermined function.

Comparison among genomes revealed a high number of proteins with conserved SPAK/OSR1 binding motifs. Indeed, 70 out of the 131 mouse proteins were also found in the human genome, indicating a 53% conservation. However, a significantly lower number of mouse and human proteins containing SPAK/OSR1 binding motifs (23 out of 131 mouse proteins) were present in the genetically more distant zebrafish genome, corresponding to 17.5% of mouse proteins. These data are summarized in Table 2 and Fig. 3.

Next, we examined the location of the SPAK/OSR1 binding motif within the entire set of 131 proteins. As seen in Fig. 4, a larger number of motifs are located at the extreme NH₂ terminus of the protein. The remaining motifs are distributed along the length of the target protein. However, proteins such as transporters and channels, which have a transmembrane core and cytosolic amino- and carboxyl termini, have the SPAK/OSR1 binding motifs preferentially located on the cytosolic tails (see Table 2). In other proteins such as kinases, which have well-defined catalytic and regulatory domains, the [S/G/V]RFx[V/I]xx[V/I/T/S]xx motifs are found exclusively in the regulatory domain. Indeed, the average location for sites found in kinases displaying SPAK/OSR1 binding motifs is $80 \pm 3.2\%$ ($n = 11$). With the exception of dynein, all SPAK/OSR1 binding motifs found in cytoskeletal proteins are also found at the extreme carboxyl-terminus (Table 2). Finally, note the opposite location of the SPAK/OSR1 binding motifs in the two regulatory subunits of protein phosphatase 2 (Table 2). These

motifs are again preferentially located at the very beginning or very end of the protein.

The identification of 137 [S/G/V]RFx[V/I]xx[V/I/T/S]xx motifs within 131 mouse proteins allow us to further examine the distribution of residues in the motif at *positions* 3 and 5–9. As shown in Fig. 5, residues at *position* 3 can be very diverse, from a mostly hydrophobic leucine to a very hydrophilic arginine. One interesting feature of this analysis of the expanded motif, is the preferential use of serine residues at *positions* 5, 6, 7, and 8.

The list of proteins that contain the SPAK/OSR1 binding motif is rather extensive despite the somewhat restrictive definition of the motif based on the 12 sequences reported in Fig. 1. However, if we were to allow additional residues at *positions* –1, 1, or 7, the total number of proteins would increase significantly. For example, we expanded our search to include the possibility of a lysine residue at *position* 1. We identified an additional 100 proteins with [S/G/V]KFx[V/I]xx[V/I/T/S]xx motifs (Table 3). To date, however, no SPAK/OSR1 interacting proteins have been identified as containing a KF_xV motif. As the motifs were identified through gene searching and not from peptide libraries, some of the identified motifs might not constitute actual SPAK/OSR1 binding sites. Therefore, whether proteins identified in this search do indeed interact with the kinases requires experimental testing. Anchoring of kinases to their substrates constitutes one important mechanism of substrate specificity (1). Thus, for any protein-protein interaction to reach biological significance, overlap of their temporal and spatial distribution is required. While SPAK/OSR1 interaction with cation-chloride cotransporters,

chloride channels, and TNF receptor has been evidenced experimentally (3–6), this study provides a list of new proteins that potentially interact with SPAK and OSR1 and suggests a significantly greater physiological role for each of these kinases.

GRANTS

This work was supported by National Institutes of Health Research Grant NS-36758.

REFERENCES

1. **Bardwell AJ, Flatauer LJ, Matsukuma K, Thorner J, Bardwell L.** A conserved docking site in MEKs mediates high-affinity binding to MAP kinases and cooperates with a scaffold protein to enhance signal transmission. *J Biol Chem* 276: 10374–10386, 2001.
2. **Delpire E, Piechotta K.** Ste20 kinases and cation-chloride cotransporters. In: *Cell Volume and Signal Transduction*, edited by Lauf PK and Adragna NC (Advances in Experimental Medicine and Biology). New York: Springer Science and Business Media 559: 43–53, 2004.
3. **Denton J, Nehrke K, Yin X, Morrison R, Strange K.** GCK-3, a newly identified Ste20 kinase, binds to and regulates the activity of a cell cycle-dependent CIC anion channel. *J Gen Physiol* 125: 113–125, 2005.
4. **Piechotta K, Garbarini NJ, England R, Delpire E.** Characterization of the interaction of the stress kinase SPAK with the Na⁺-K⁺-2Cl[–] cotransporter in the nervous system: Evidence for a scaffolding role of the kinase. *J Biol Chem* 278: 52848–52856, 2003.
5. **Piechotta K, Lu J, Delpire E.** Cation-chloride cotransporters interact with the stress-related kinases SPAK and OSR1. *J Biol Chem* 277: 50812–50819, 2002.
6. **Polek TC, Talpaz M, Spivak-Kroizman T.** The TNF receptor, RELT, binds SPAK and uses it to mediate p38 and JNK activation. *Biochem Biophys Res Commun* 343: 125–134, 2006.
7. **Waterston RH, Mouse Genome Sequencing Consortium.** Initial sequencing and comparative analysis of the mouse genome. *Nature* 420: 520–562, 2002.