

<b>Manuscript Number:</b>	GICS-D-17-01165R2	
<b>Full Title:</b>	A chelicerate-specific burst of nonclassical Dscam diversity	
<b>Article Type:</b>	Research article	
<b>Section/Category:</b>	Comparative and evolutionary genomics	
<b>Funding Information:</b>	the National Natural Science Foundation of China (31630089, 31430050, 91740104)	Mr Yongfeng Jin
<b>Abstract:</b>	<p><b>Abstract</b></p> <p><b>Background</b></p> <p>The immunoglobulin (Ig) superfamily receptor Down syndrome cell adhesion molecule (Dscam) gene can generate tens of thousands of isoforms via alternative splicing, which is essential for both nervous and immune systems in insects. However, further information is required to develop a comprehensive view of Dscam diversification across the broad spectrum of Chelicerata clades, a basal branch of arthropods and the second largest group of terrestrial animals.</p> <p><b>Results</b></p> <p>In this study, a genome-wide comprehensive analysis of Dscam genes across Chelicerata species revealed a burst of nonclassical Dscams, categorised into four types: mDscam, sDscam<math>\alpha</math>, sDscam<math>\beta</math>, and sDscam<math>\gamma</math> based on their size and structure. Although the mDscam gene class includes the highest number of Dscam genes, the sDscam genes utilise alternative promoters to expand protein diversity. Furthermore, we indicated that the 5' cassette duplicate is inversely correlated with the sDscam gene duplicate. We showed differential and sDscam-biased expression of nonclassical Dscam isoforms. Thus, the Dscam isoform repertoire across Chelicerata is entirely dominated by the number and expression levels of nonclassical Dscams. Taken together, these data show that Chelicerata evolved a large conserved and lineage-specific repertoire of nonclassical Dscams.</p> <p><b>Conclusions</b></p> <p>This study showed that arthropods have a large diversified Chelicerata-specific repertoire of nonclassical Dscam isoforms, which are structurally and mechanistically distinct from those of insects. These findings provide a global framework for the evolution of Dscam diversity in arthropods and offer mechanistic insights into the diversification of the clade-specific Ig superfamily repertoire.</p>	
<b>Corresponding Author:</b>	Yongfeng Jin, PH.D Zhejiang University College of Life Sciences HangZhou, Zhejiang CHINA	
<b>Corresponding Author Secondary Information:</b>		
<b>Corresponding Author's Institution:</b>	Zhejiang University College of Life Sciences	
<b>Corresponding Author's Secondary Institution:</b>		
<b>First Author:</b>	Guozheng Cao	
<b>First Author Secondary Information:</b>		
<b>Order of Authors:</b>	Guozheng Cao	
	Yang Shi	
	Jian Zhang	
	Hongru Ma	
	Shouqing Hou	

	Haiyang Dong
	Weiling Hong
	Shuo Chen
	Hao Li
	Yandan Wu
	Pengjuan Guo
	Xu Shao
	Bingbing Xu
	Feng Shi
	Yijun Meng
	Yongfeng Jin, PH.D
<b>Order of Authors Secondary Information:</b>	
<b>Response to Reviewers:</b>	The authors' response letter has been included as a supplementary file

[Click here to view linked References](#)

# A chelicerate-specific burst of nonclassical Dscam diversity

Guozheng Cao<sup>1§</sup>, Yang Shi<sup>1§</sup>, Jian Zhang<sup>1</sup>, Hongru Ma<sup>1</sup>, Shouqing Hou<sup>1</sup>, Haiyang Dong<sup>1</sup>,  
Weiling Hong<sup>1</sup>, Shuo Chen<sup>1</sup>, Hao Li<sup>1</sup>, Yandan Wu<sup>1</sup>, Pengjuan Guo<sup>1</sup>, Xu Shao<sup>1</sup>, Bingbing Xu<sup>1</sup>,  
Feng Shi<sup>1</sup>, Yijun Meng<sup>2</sup>, Yongfeng Jin<sup>1\*</sup>

<sup>1</sup>Institute of Biochemistry, College of Life Sciences, Zhejiang University, Hangzhou, Zhejiang,  
ZJ310058, <sup>2</sup>College of Life and Environmental Sciences; Hangzhou Normal University;  
Hangzhou, Zhejiang, ZJ310018, P. R. of China,

\* Correspondence should be addressed to Yongfeng Jin: 0086-571-88206479(Tel);  
0086-571-88206478(Fax); [jinyf@zju.edu.cn](mailto:jinyf@zju.edu.cn) (e-mail).

§ These authors contributed equally to this work

Guozheng Cao: [15275247155@163.com](mailto:15275247155@163.com); Yang Shi: [18967116958@163.com](mailto:18967116958@163.com); Shouqing Hou:  
[17816860428@163.com](mailto:17816860428@163.com); Yijun Meng: [mengyijun@zju.edu.cn](mailto:mengyijun@zju.edu.cn); Hao Li: [aholdlee@163.com](mailto:aholdlee@163.com);  
Jian Zhang: [zhangjian2301@163.com](mailto:zhangjian2301@163.com); Hongru Ma: [13588702934@163.com](mailto:13588702934@163.com); Haiyang Dong:  
[21507070@zju.edu.cn](mailto:21507070@zju.edu.cn); Weiling Hong: [hongweiling1989@126.com](mailto:hongweiling1989@126.com); Chenshuo:  
[chenshuo124@126.com](mailto:chenshuo124@126.com); Yandan Wu: [2601463242@qq.com](mailto:2601463242@qq.com); Pengjuan Guo:  
[guopengjuan0613@163.com](mailto:guopengjuan0613@163.com); Xu Shao: [therenothing@foxmail.com](mailto:therenothing@foxmail.com); Bingbing Xu:  
[xubb1987@126.com](mailto:xubb1987@126.com); Feng Shi: [shifeng@zju.edu.cn](mailto:shifeng@zju.edu.cn); Yongfeng Jin: [jinyf@zju.edu.cn](mailto:jinyf@zju.edu.cn);

Running title: A burst of nonclassical Dscam diversity in chelicerate

Keywords: Dscam, Gene duplication, Exon duplication, Alternative splicing, Alternative  
promoter, Chelicerate

# Abstract

## Background

The immunoglobulin (Ig) superfamily receptor Down syndrome cell adhesion molecule (*Dscam*) gene can generate tens of thousands of isoforms via alternative splicing, which is essential for both nervous and immune systems in insects. However, further information is required to develop a comprehensive view of *Dscam* diversification across the broad spectrum of Chelicerata clades, a basal branch of arthropods and the second largest group of terrestrial animals.

## Results

In this study, a genome-wide comprehensive analysis of *Dscam* genes across Chelicerata species revealed a burst of nonclassical *Dscams*, categorised into four types—*mDscam*, *sDscam* $\alpha$ , *sDscam* $\beta$ , and *sDscam* $\gamma$ —based on their size and structure. Although the *mDscam* gene class includes the highest number of *Dscam* genes, the *sDscam* genes utilise alternative promoters to expand protein diversity. Furthermore, we indicated that the 5' cassette duplicate is inversely correlated with the *sDscam* gene duplicate. We showed differential and *sDscam*-biased expression of nonclassical *Dscam* isoforms. Thus, the *Dscam* isoform repertoire across Chelicerata is entirely dominated by the number and expression levels of nonclassical *Dscams*. Taken together, these data show that Chelicerata evolved a large conserved and lineage-specific repertoire of nonclassical *Dscams*.

## Conclusions

This study showed that arthropods have a large diversified Chelicerata-specific repertoire of nonclassical Dscam isoforms, which are structurally and mechanistically distinct from those of insects. These findings provide a global framework for the evolution of Dscam diversity in arthropods and offer mechanistic insights into the diversification of the clade-specific Ig superfamily repertoire.

## Keywords

Dscam, Gene duplication, Exon duplication, Alternative splicing, Alternative promoter, Chelicerate

## Background

Alternative splicing plays an important role in the generation of proteomic diversity and genomic evolution in metazoans [1,2]. For example, the Down syndrome cell adhesion molecule (*Dscam*) gene in *Drosophila melanogaster* has the potential to generate 38,016 distinct mRNA and protein isoforms via mutually exclusive alternative splicing [3]. In this *Dscam* gene structure, 95 alternatively spliced exons are organised into exon 4, 6, 9, and 17 clusters that contain 12, 48, 33, and 2 copies, respectively. Forward genetic screening and biochemical approaches have shown that this extensive diversity of Dscam encoded by a single locus is required for both nervous and immune systems. Dscam isoform diversity plays an important role in neuronal wiring and self-recognition, and the extensive diversity of Dscam isoforms has been shown to confer specificity for antigen recognition [4–13].

1 Vertebrate Dscams, which lack the striking diversity of their insect counterparts, are  
2  
3 involved mainly in the developmental processes of the nervous system [14]. Interestingly,  
4  
5 cadherin superfamily members, protocadherins (Pcdhs), might serve an analogous function in  
6  
7  
8  
9 vertebrates [15–17]. Human *Pcdh* genes are arranged tandemly in three groups called *Pcdh $\alpha$* ,  
10  
11 *Pcdh $\beta$* , and *Pcdh $\gamma$* , with 14, 22, and 22 repeats in their respective 5' variable regions [18].  
12  
13 Unlike *Drosophila Dscam1*, these *Pcdhs* utilise alternative promoters to generate isoform  
14  
15 diversity [19,20]. For both Dscams and Pcdhs, isoform expression appears largely stochastic  
16  
17 and combinatorial rather than determinative, endowing a unique cell-surface identity for each  
18  
19  
20 neuron [15,16,21,22]. Both molecules exhibit isoform-specific homophilic binding, and they  
21  
22  
23 may have similar roles in the nervous system. Interestingly, *Drosophila* lacks the counterparts  
24  
25  
26 of vertebrate *Pcdhs*. Thus, these two phyla appear to have independently evolved similar  
27  
28  
29 molecular strategies for comparable roles by recruiting various molecules from different  
30  
31  
32 protein families [15].  
33  
34  
35  
36  
37

38 As Pcdh isoform diversity is restricted to vertebrates, Dscam isoform diversity has been  
39  
40 considered unique to arthropods [23,24]. However, compared with the well-known  
41  
42 phylogenetic distribution of vertebrate *Pcdhs*, *Dscam* diversification has not been studied in  
43  
44 detail in arthropods, particularly in Chelicerata members, which represent a basal branch of  
45  
46 arthropods and form the second largest group of terrestrial animals. It is generally recognised  
47  
48  
49 that Dscam protein structure is conserved across bilaterians, containing 10 immunoglobulin  
50  
51 (Ig) domains and six fibronectin III (FNIII) repeats, with the tenth Ig domain located between  
52  
53 FNIII 4 and FNIII 5 [14,15,17]. Recently, structural variants of Dscam have been found in  
54  
55  
56 *Ixodes scapularis* [24]. We identified shortened *Dscam* genes with tandemly arrayed 5'  
57  
58  
59  
60  
61  
62  
63  
64  
65

cassettes in several Chelicerata species, similar to vertebrate clustered *Pcdhs* [25]. These results suggest that the structural and expansion pathways of *Dscam* might be diversified in Chelicerata species. However, no comprehensive view exists of *Dscam* diversification across the broad spectrum of Chelicerata. The rapidly increasing availability of genomic information regarding Chelicerata, particularly key clades such as the Xiphosuran horseshoe crab, will increase our understanding of *Dscam* diversification.

In this study, we performed a genome-wide comprehensive analysis of *Dscam* genes in chelicerates, which diverged from other Arthropod lineages ~500 million years ago. The identified *Dscam* genes could be grouped into one classical (*LDscam*) and four nonclassical (*mDscam*, *sDscama*, *sDscamβ*, and *sDscamγ*) types based on their size and structure. Although the *mDscam* gene class includes the highest number of *Dscam* genes, the *sDscam* genes utilise alternative promoters to expand protein diversity. Collectively, these results demonstrate that Chelicerata specifically evolved different organisation and mechanisms that generated a diverse lineage-specific repertoire of *Dscam* isoforms. These findings highlight the rich Chelicerata-specific diversification of *Dscam* genes, and provide a global framework for the evolution of *Dscam* diversity in arthropods and bilaterians.

## Results

### Genome-wide identification of *Dscam* genes across Chelicerata species

To generate a global blueprint for *Dscam* diversity in Chelicerata, we performed a genome-wide analysis of *Dscam* homologues in representative species from each of the major

clades. We examined one species of the order Merostomata (*Limulus polyphemus*) and five species representing five major clades of the order Arachnida, including two Araneae (*Stegodyphus mimosarum* and *Parasteatoda tepidariorum*), one Scorpiones (*Mesobuthus martensii*), one Mesostigmatan (*Metaseiulus occidentalis*), and one Ixodidan (*Ixodes scapularis*) (Additional file 1: Table S1). These organisms constitute some of the major taxonomic groups of the Chelicerata subphylum that last shared a common ancestor ~500 million years ago [26], with long-term resolution (comparing the gene organisation among the five Arachnida clades and Merostomata). Using cross-species comparisons with RNA-sequencing (RNA-seq) analyses, we identified 198 *Dscam* genes in six representatives of the Chelicerata species, 161 of which were novel or corrected (Fig. 1b; Additional file 2: Table S2). Our results indicated that all extant chelicerates display marked expansion of the *Dscam* gene family.

Phylogenetic analysis indicated that all of the *Dscam* proteins could be clustered into three groups: canonical *Dscams* and two groups of shortened nonclassical proteins (Fig. 2). The three groups differ in size, and are hereafter referred to as LDscam, mDscam, and sDscam, respectively (Fig. 1a). The sDscams lack the canonical Ig1–6,10 and FNIII 3–4,6 domains present in classical *Dscam*. The genes encoding these sDscams can be subdivided into three types—*sDscam* $\alpha$ , *sDscam* $\beta$ , and *sDscam* $\gamma$ —based on differing structures in the 5' regions (Fig. 1a). *sDscam* $\alpha$  and *sDscam* $\beta$  are characterised by clustered cassette repeats in the 5' regions, encoding one and two Ig domains, respectively. *sDscam* $\gamma$  shares similar domains with *sDscam* $\alpha$  and *sDscam* $\beta$ , albeit without tandemly arrayed cassettes in the 5' regions. Another type of shortened *Dscam* lacks the FNIII and Ig domains present in the C-terminal region of



classical Dscam. As the sizes of these Dscams fall between those of canonical Dscam (LDscam) and sDscam, we designated these intermediate shortened *Dscam* genes as *mDscams*. As shown in Fig. 1b, nonclassical Dscams dominate the isoform repertoire in all Chelicerata species investigated.

### **Nonclassical *Dscam* genes are Chelicerata-conserved and restricted**

This study revealed that nonclassical *Dscams* are conserved across Chelicerata species and that all types of *sDscams* (*sDscam* $\alpha$ , *sDscam* $\beta$ , and *sDscam* $\gamma$ ) are largely present in representative species from the Arachnida order and Merostomata class investigated in this study (Fig. 1b). Our results suggest that these types of nonclassical *Dscams* are ancient, existing before the split of Arachnida and Merostomata. The *Dscam* gene family members diverged markedly across various species. The higher gene number in scorpions, spiders, and horseshoe crabs is consistent with an additional event of whole or large-scale genomic duplications [27]. Although *mDscam* genes are present in greater numbers than *LDscams* and *sDscams*, the *sDscam* genes can generate up to a hundred isoforms through a combination of alternating promoters and alternative splicing. Thus, based on the isoform number, *sDscams* are represented more than *LDscams* and *mDscams* in each Chelicerata. However, no such nonclassical *Dscams* have been identified among the *Dscam* genes from the Mandibulata species of insect, Crustacea, or Myriapoda classes, or in any non-arthropod species. This observation suggests that they arose after the radiation of Mandibulata and Chelicerata during Arthropoda evolution. Thus, we conclude that the nonclassical *Dscams* are largely conserved and restricted to Chelicerata.

## Global analysis of *Dscam* relationships over the Chelicerata phylogeny

To trace the evolutionary history of nonclassical *Dscam* genes, we first performed multiple sequence alignments of all *Dscams* from six representative species. Comparative analysis of *Dscam* sequences from Chelicerata and outgroup species revealed three major clades, which represent three groups of *LDscam*, *mDscam*, and *sDscam* (Fig. 2). Based on the phylogenies of individual *Dscam* types (Additional file 3: Figs. S1, 2), although our analyses of the *Dscam* genes did not permit a full ancestral reconstruction, several conclusions could be reached. First, these data indicate that at least seven *mDscams* and five *sDscams* were present in the Chelicerata ancestor before the split of Arachnida and Merostomata (Additional file 3: Figs. S1, 2). Second, three types of *LDscams*, *mDscams*, and *sDscams* were inconsistently expanded. *mDscam* genes have undergone massive duplications during Chelicerata evolution, while *LDscam* genes have undergone few or limited duplications (Fig. 2; Additional file 3: Figs. S1, 2). For example, the divergence of the Araneae ancestor into the *P. tepidariorum* and *S. mimosarum* ancestors involved 10 *P. tepidariorum*-specific duplications of *mDscam* (Additional file 3: Fig. S1). This gene expansion process is ongoing, as demonstrated by recent duplications in the Araneae species.

Phylogenetic analysis revealed that *sDscams* could be clustered into four clades (clades A–D; Additional file 3: Fig. S2). Clade D exclusively included *sDscam*βs from all species investigated, suggesting this *sDscam*β is ancient and arose in the Chelicerata ancestor. Clade B consisted of conserved *sDscamas* and species-specific *sDscamy*, suggesting that this *sDscama* arose from the Chelicerata ancestor. Interestingly, clade C included species-specific

*sDscamβs* and *sDscamas*, in addition to *sDscamys*. We speculate that an *sDscamy* ancestor evolved differentially into *sDscamβs* or *sDscamas* during Chelicerata divergence. These results indicated that *sDscamβs* and *sDscamas* might have multiple independent origins.

### **A lineage-specific burst of 5' clustered cassettes in *sDscam* genes**

To produce an overview of the evolutionary relationships among the variable cassettes, we generated heatmaps for *sDscama* and *sDscamβ* to show the relative similarities of each cassette repeat to other variable repeats both within and between species. For these analyses, we selected one representative species from each of the major orders investigated: *M. occidentalis*, *S. mimosarum*, *M. martensii*, *I. scapularis*, and *L. polyphemus*. Analysis of the heatmaps of the tandemly arrayed 5' cassettes in the *sDscamas* and *sDscamβs* revealed little evidence of conserved orthologous pairs of repeats between species (Fig. 3a, b). Instead, when striking similarities were found between repeats in each species, they typically involved large blocks of highly similar cassettes within each gene of one species. For example, a massive block of 61, 36, and 18 cassettes expanded specifically in *sDscama* of *M. occidentalis*, and in *sDscama1* and *sDscama2* of *S. mimosarum*, respectively (Fig. 3c). Similarly, blocks of 21 and 33 cassettes in the *M. occidentalis* *sDscamβ1* and *sDscamβ2*, respectively, were highly similar to one another (Fig. 3b, d). Phylogenetic analysis of duplicated cassettes from two closely related spiders (*S. mimosarum* and *P. tepidariorum*) indicated that 5' cassette duplication largely occurs in a lineage-specific manner (Additional file 3: Fig. S3). Furthermore, duplicated cassettes showed the tendency to be located adjacent to one another within a gene. Overall, these data indicate that the main expansions of

ancestral cassettes of *sDscama* and *sDscamβs* occurred independently in each lineage.

The heatmap patterns differed considerably between *sDscamas* and *sDscamβs* among the various species. The *sDscama* heatmap revealed large blocks of cassette duplications specific to each clade, consistent with phylogenetic analyses of duplicated cassettes across Chelicerata species (Fig. 3c). In contrast, with the exceptions of *M. occidentalis*, most *sDscamβs* contain a number of common blocks that span multiple clades, as represented in green shading in Figure 4b. Notably, duplications of 5' cassettes in *sDscamas* and *sDscamβs* are almost exclusively specific to *M. occidentalis*, showing fewer similarities with other cassettes from other species, as represented by blue blocks in the regions of the heatmap that compare two species (Fig. 3a, b). Such high species specificity of the 5' variable cassettes suggests that *sDscamas* and *sDscamβs* have undergone rapid expansions and divergence. Furthermore, comparison of 5' cassette- and gene-based clustering clearly indicated that 5' cassette duplication occurs on a much faster timescale than gene duplication (Additional file 3: Fig. S4). This multi-layer expansion of sDscam diversity might help to increase the efficiency and flexibility of spatiotemporal regulation.

### **5' cassette duplicate inversely correlated with *sDscam* gene duplicates**

As gene duplications and 5' cassette tandem duplications increased isoform diversity through the emergence of additional genomic copies, we next investigated whether or how they are related to each other as evolutionary mechanisms. We found that the number of 5' cassette tandem duplicates correlated inversely with that of gene duplicates in the *sDscama* and *sDscamβ* subfamily (Fig. 4a, b). Single *sDscams* (singletons) are likely to contain

substantially more tandem cassettes in the 5' variable region. For example, up to 62 and 40 copies reside in the 5' variable region of *sDscama* in *M. occidentalis* and *M. martensii*, respectively. Furthermore, the *sDscama* subfamily (1–3 members) contains two- to four-fold more 5' cassette duplicates per gene, which is similar to members in the larger *sDscamβ* subfamily (2–7 members) (Fig. 4c); we found the reverse trend in gene duplicates (Fig. 4d). Importantly, the total number of Dscam isoforms is roughly similar among various Chelicerata species (Fig. 1b). Therefore, this correlation may reflect compensatory evolution between alternative promoter and gene duplication, analogous to the inverse correlation of alternative splicing with gene duplication [28]. These results imply that the inverse correlation is not simply an inherent inclination, but is instead fulfilling the complementary demand for expanding Dscam isoforms via distinct evolutionary mechanisms.

To further examine the evolutionary forces underlying the inverse correlation, we carried out a detailed comparison of duplication scenarios in *sDscama* and *Dscamβ* genes from various species. Interestingly, the number of 5' cassette tandem duplicates per gene correlated inversely with the size of the duplicate blocks in *M. occidentalis*, *I. scapularis*, *M. martensii*, and *S. mimosarum* (Fig. 4e–h). These correlations largely fitted to a power law. *L. polyphemus* did not exhibit as strong a correlation as the other species (Fig. 4i), possibly due to the incomplete annotation of the *Dscam* genes. Furthermore, the numbers of 5' cassette duplicates per *sDscama* were roughly 1–3-fold higher than those in *sDscamβ*s, with the duplicate blocks containing two additional exons than those in *sDscamas* (Fig. 4j). This inverse relationship might reflect an inherent property of the species, *i.e.*, the genome size (Fig. 4k, l). For example, *I. scapularis*, with a large genome size estimated at 2,100 Mb [29],

1 contained more *sDscam* gene duplicates, but fewer 5' cassette tandem duplicates in each gene.  
2  
3 In contrast, *M. occidentalis*, with a small genome size (~152 Mb), had fewer *sDscam* gene  
4  
5 duplicates, but more 5' cassette tandem duplicates per gene. In particular, phylogenetic  
6  
7 analysis indicated loss of cassette-within introns occurred independently in *sDscamβ1* and  
8  
9 *sDscamβ2* of *M. occidentalis*, whereas no substantial intron loss occurred in the constant  
10  
11 region (Fig. 4m). This result led us to speculate that intron loss caused the decreased repeat  
12  
13 block size to facilitate greater duplication in the 5' variable region. Conversely, this  
14  
15 species-specific intron loss might be driven by the selection pressure of greater cassette  
16  
17 duplication. These results suggest that gene duplications and 5' cassette tandem duplication  
18  
19 are not selected to expand independently of each other during Chelicerata evolution.  
20  
21  
22  
23  
24  
25  
26  
27  
28

## 29 **Differential and biased expression of nonclassical Dscam isoforms**

30  
31  
32  
33 To further characterise the expression of nonclassical Dscam isoforms, we employed  
34  
35 RNA-seq data to analyse the expression of the nonclassical *Dscams* in various tissues of *M.*  
36  
37 *martensii*. Similar to *sDscama* and *sDscamβ* [25], as well as three classical *Dscams*  
38  
39 (*LDscams1–3*), 14 of 17 *mDscams* and 3 of 4 *sDscamys* were expressed at markedly higher  
40  
41 levels in the cephalothorax than in other tissues (Fig. 5a; Additional file 3: Fig. S5). This  
42  
43 expression pattern is largely coincident with high expression of classical Dscams in the  
44  
45 nervous system of vertebrates and insects [5,30–32]. In contrast, *LDscams*, *mDscams*, and  
46  
47 *sDscams* showed low expression in muscles. Notably, *mDscam8* was specifically expressed at  
48  
49 maximum levels in haemocytes, whereas *mDscam10* and *mDscam15* were highly expressed  
50  
51 in poison glands (Fig. 5a). It will be interesting to see whether nonclassical *Dscams* play a  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

potential role in Scorpione immunity, similar to that of *Dscam1* isoforms in the *Drosophila* immune system [5]. These results suggest that nonclassical *Dscam* expression is regulated spatially and temporally.

The wealth of published RNA-seq data from a set of 20 embryo-derived cell lines allows the comparison of *Dscam* expression in various cell lines of *I. scapularis* [29] (Additional file 4: Table S3). Expression profiling revealed that *Dscam* genes were differentially expressed in various cell lines (Fig. 5b). Interestingly, many *Dscams* were preferentially and specifically expressed in certain lines. For example, high expression of *mDscam14* was observed in strain IDE2, but was almost undetectable in strain IDE8 and ISE18. In contrast, *sDscamγ4* showed robust expression in strain IDE8 and ISE18, but low level expression in strain IDE12; *sDscamβ2* showed high expression in strain ISE18. Likewise, the 5' variable exons of *sDscama* and *sDscamβ* genes exhibited differential expression in various lineages (Additional file 3: Fig. S6). *sDscama2.14* was specifically expressed at maximum levels in strain ISE6, while *sDscamβ2.12* and *sDscamβ6.6* were expressed at maximum levels in strain ISE18 and IDE12, respectively. Collectively, these results revealed lineage-specific expression signatures of *Dscam* isoforms.

We observed a dramatic bias: *sDscams* were largely expressed at higher levels than *LDscams* and *mDscams* in various tissues of *M. martensii* (Fig. 6a), and in different *I. scapularis* lineages (Fig. 6b). A similar bias was observed in *S. mimosarum*, *P. tepidariorum*, *M. occidentalis*, and *L. polyphemus* (Fig. 6c). These data indicate that *sDscam*-biased expression patterns are evolutionally conserved across Chelicerata. It is possible that the high expression

of *sDscams* might be due to the use of multiple promoters, leading to more transcripts. Taken together, both the number and expression level of nonclassical *Dscams* dominate exclusively in the Dscam isoform repertoire across Chelicerata.

## Discussion

### The evolutionary landscape of Dscam diversity in arthropods

This study identified large and diverse nonclassical Dscam repertoires in chelicerates, most of which are newly annotated or corrected. The results reported here extend our previous insights into Dscam diversity during arthropod evolution [23–25]. Firstly, the nonclassical *Dscam* genes can be classified into four types (*mDscam*, *sDscam $\alpha$* , *sDscam $\beta$* , and *sDscam $\gamma$* ) based on their size and structure (Fig. 1). Moreover, the nonclassical *Dscam* genes are conserved and restricted to Chelicerata, suggesting that Chelicerata uniquely evolved a large lineage-specific repertoire of nonclassical Dscam isoforms. Phylogenetic analysis and comparison of Dscam structures revealed that the *Dscam* ancestor underwent multiple shortening events during chelicerate evolution, leading to the loss of protein domains to varying degrees [25]. Thus, the *Dscam* gene number has undergone massive expansion, and also the structure of *Dscam* genes became highly diversified during Chelicerata evolution. These findings, together with those of others [23–25], provide a global framework for the evolution of Dscam diversity in arthropods.

These Dscam diversifications suggest that their binding mechanisms are potentially different from those of classical Dscams. The crystal structural analysis reveals that the first



four Ig domains of *Drosophila* Dscam1 are responsible for the formation of the horseshoe configuration. Additional modelling studies illuminated the molecular basis of the isoform-specific homophilic binding specificity [15]. As they contain the same eight N-terminal Ig domains as *Drosophila* Dscam1, it is conceivable that the chelicerate mDscams studied in this research could form a horseshoe configuration and interact via a similar mechanism. However, this horseshoe structure does not form in chelicerate sDscams with only three Ig domains. Therefore, we speculate that chelicerate sDscams exhibit a different mode of isoform-specific homophilic binding than *Drosophila* Dscam1.

Chelicerata seem to have generated far fewer Dscam isoforms than insects. As estimated by the number of Ig7 or orthologues, the number of Dscam isoforms is in the range of ~100–200 across the Chelicerata species investigated (Fig. 1), approximately two to three-fold lower than that in insects. Recent studies indicate that clustered mammalian Pcdhs could expand the binding specificity repertoire via *cis*-multimers. For 22  $\gamma$ -Pcdhs, the diversity of adhesive interfaces could be on the order of  $10^5$  through *cis*-tetramerisation coupled with haemophilic *trans* interactions [33,34]. Given the striking organisational resemblance between the Chelicerata clustered *sDscams* and mammalian *Pcdhs*, it is attractive to speculate that Chelicerata sDscams could function via *cis*-multimers. If so, the diversity of adhesive interfaces mediated by Chelicerata sDscams will be much higher, as there are many more Chelicerata-clustered *sDscams* (90–130) than mammalian-clustered *Pcdhs* (50–60).

## Dscams versus Protocadherins

1 Dscams and Pcdhs belong to large established families of cell adhesion molecules: Dscams  
2  
3 belong to the Ig superfamily and Pcdhs belong to the cadherin superfamily [35]. The clustered  
4  
5 Pcdh diversity is confined within the clade of jawed vertebrates, which is considered as a  
6  
7 chordate innovation [36], whereas extensive Dscam diversity is unique to arthropods. In the  
8  
9 latter case, arthropods diversify using two mechanisms to generate Dscam isoform diversity:  
10  
11  
12 Mandibulata *Dscam* genes employ exclusive splicing of internal exon clusters to generate  
13  
14 distinct isoforms [3,24] and chelicerate *Dscams* utilise alternative promoters in the 5' variable  
15  
16 region [25]. However, neither cadherin nor protocadherin genes appear to have internal  
17  
18 tandem exon arrays in a manner similar to Mandibulata *Dscams*.  
19  
20  
21  
22  
23  
24  
25

26 Both clustered chelicerate *Dscams* and vertebrate *Pcdhs* are organised in a tandem array in  
27  
28 the 5' variable region (Fig. 7a, b). Curiously, both 5' clustered *sDscams* and *Pcdhs* appeared to  
29  
30 originate via an analogous evolutionary pathway (Fig. 7c, d), which was involved in the  
31  
32 shortening and expansion of *Dscam* and *cadherin* ancestors [25,37]. In both genes, each  
33  
34 variable repeat is preceded by a promoter, and differential expression occurs via combining  
35  
36 alternate promoter choice with alternative splicing [19,20,25]. Moreover, 5' clustered *Dscams*  
37  
38 and *Pcdhs* contain similar structural composition encoding six extracellular domains, a single  
39  
40 transmembrane (TM) region, and a cytoplasmic domain. Interestingly, three-dimensional  
41  
42 protein structure modelling revealed a similar  $\beta$  sandwich structure between the first domain  
43  
44 of Pcdh and sDscam (Fig. 7a, b). Finally, we showed that clustered *sDscams* encode proteins  
45  
46 exhibiting isoform-specific homophilic binding in a manner similar to *Pcdhs*. Despite their  
47  
48 overall similarities, the structural properties of *Pcdh* and *sDscam* genes differ in at least two  
49  
50 major aspects. First, inconsistent with the conserved arrangement of a single genome locus  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1 containing three tandem *Pcdh* gene clusters across vertebrates, the *sDscams* are largely  
2  
3 dispersed or partially clustered in the chelicerate genome. As the transcription of *Pcdh* gene  
4  
5 clusters is closely linked and mediated by long-range chromatin-looping interactions [38], the  
6  
7 dispersed distribution of *sDscam* gene clusters might reflect a mode of transcription  
8  
9 regulation distinct from that of *Pcdh* genes. A second major difference concerns the structure  
10  
11 of the variable region. *Pcdh* variable exons encode the entire ectodomain composed of six  
12  
13 extracellular cadherin domains (EC1–EC6), a single TM region, and a short cytoplasmic  
14  
15 extension, whereas the variable cassettes of *sDscams* encode the partial ectodomain of one or  
16  
17 two Ig domains at the N-terminus. Future studies are needed to investigate the role of variable  
18  
19 region structures in the subcellular distribution, isoform-specific binding, and multimer  
20  
21 formation.  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31

32 Combining experimental evidence with the complementary phylogenetic distribution of  
33  
34 Dscam diversity in arthropods and Pcdh diversity in vertebrates, it is tempting to suggest that  
35  
36 both may have similar roles in the nervous system [15–17,21,22]. These two phyla seem to  
37  
38 employ a similar strategy for self/non-self discrimination by recruiting various molecules of  
39  
40 different protein families [15]. Nevertheless, there is a wide evolutionary gap between  
41  
42 arthropods and vertebrates, as they share a common ancestor more than 500 million years ago.  
43  
44 It will be informative to explore the molecules or mechanisms of species within this  
45  
46 phylogenetic gap (i.e., *Branchiostoma floridae*), which are thought to lack both clustered  
47  
48 *Pcdh* and *Dscam* genes, and that have evolved to endow cells with distinct molecular  
49  
50 identities and highly diverse recognition selectivity strategies.  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

## Diversification of Ig superfamily protein repertoires

Although tandemly arrayed Ig repeats are found frequently in the genome, animals achieve protein diversity from a single locus via a variety of mechanisms (Additional file 3: Fig. S7).

In higher vertebrates, the great diversity of antigen-specific receptors of the adaptive immune system can be achieved through somatic gene rearrangement and clonal selection at the DNA transcriptional level, which is known as the V(D)J mechanism [39–42]. However, insect *Dscaml*s utilise mutually exclusive splicing to generate an extensive repertoire of thousands of Ig-superfamily protein isoforms [3,5]. In contrast, our studies indicate that Chelicerata *sDscams* employ alternative promoters to generate substantial numbers of isoforms [25]. These molecular processes are mutually exclusive among three distinct clades, functioning to create extensive molecule diversity. Metazoans might have evolved different ways for extensive Ig-superfamily proteins to enable immune defence and other functions.

## Conclusions

In this study, we identified large and lineage-specific nonclassical Dscam repertoires in chelicerates. Nonclassical *Dscams* are conserved and restricted to Chelicerata, and have been classified into four types based on their size and structure. These results demonstrate that arthropods specifically diversify a large Chelicerata-specific repertoire of nonclassical Dscam isoforms. The Dscam isoform repertoire across Chelicerata is dominated exclusively by the number and expression levels of nonclassical *Dscams*. This genome-wide identification and classification study of *Dscam* genes provides the global framework of the evolution of Dscam diversity in arthropods, and provides mechanistic insights into the diversification of the

species-specific Ig superfamily repertoire.

## METHODS

### Data availability and RNA-seq data analysis

We investigated the following representative Chelicerate species: Mesostigmatan *M. occidentalis*, Trombidiformes *I. scapularis* [29], two Araneae *S. mimosarum* and *P. tepidariorum* [43], two Scorpiones *M. martensii* [44], and Merostomatan *L. polyphemus*. The sources of the Chelicerata genome sequences used in this study are shown in Table S1 (Additional file 1). To validate the *Dscam* candidates, we selected 125 publically available RNA-seq datasets corresponding to various developmental stages, tissues, organs, and cell lines across six chelicerate species (Additional file 4: Table S3). All of the raw RNA-seq datasets were subject to pre-treatment, including adapter trimming and low-quality read removal using the FASTX-Toolkit ([http://hannonlab.cshl.edu/fastx\\_toolkit/index.html](http://hannonlab.cshl.edu/fastx_toolkit/index.html)). Then, for each sequencing dataset, the count of each RNA-seq read was normalized to reads per million (RPM), thus enabling cross-sample comparison of the *Dscam* expression levels. Specifically, to calculate the RPM of one RNA-seq read, the raw count of the read was divided by the total raw count of the RNA-seq dataset and multiplied by  $10^6$ . The treated reads were mapped onto the transcripts using Bowtie 2 [45]. From these mappings, we were able to calculate the expression levels in reads per kilobase of transcript per million mapped reads (RPKM).

### Availability of genome and RNA-seq data

We investigated the following representative species of Chelicerate: Mesostigmatan *M. occidentalis*, Trombidiformes *I. scapularis* [29], two Araneae *S. mimosarum* and *P. tepidariorum* [43], two Scorpiones *M. martensii* [44], and Merostomatan *L. polyphemus*. The sources of the Chelicerata genome sequences used in this study are shown in Table S1 (Additional file 1). For *Dscam* candidate validation, we selected 125 publically available RNA-seq data corresponding to various developmental stages, tissues, and organs, and cell lines across six chelicerate species (Additional file 4: Table S3).

### Annotation and identification of *Dscam* genes

The sequences of the *Dscam* homologues were annotated through cross-species BLAST searches using the available annotated *Dscam* sequences (<http://blast.ncbi.nlm.nih.gov/Blast.cgi>). These *Dscam* candidate homologues were validated further using publically available transcriptome and RNA-seq datasets. All *Dscam* candidates were confirmed by phylogenetic analysis (<http://www.ebi.ac.uk/clustalw/index.html>) and then analysed by classifying and predicting protein domains with InterPro [46] (<http://www.ebi.ac.uk/interpro/>) and PROSITE [47] (<http://prosite.expasy.org/prosite.html>). The *Dscam* genes identified in representative species of Chelicerate are listed in Table S2 (Additional file 2).

### Phylogenetic analysis

The *Dscam* sequences were aligned across species using the Clustal W2 software package (<http://www.ebi.ac.uk/clustalw/index.html>) [48]. The coding sequences of the variable region

1 were translated, and the resulting polypeptides were aligned. The nucleotide sequences of  
2  
3 each 5' variable cassette of *sDscam* were translated into amino acid sequences and aligned.  
4  
5  
6 The genetic distances for each gene were estimated using MEGA 7.0 software [49]. We used  
7  
8  
9 maximum likelihood (ML) methods, using MEGA [49], to build the phylogenies. For the ML  
10  
11  
12 analysis, we ran MEGA with at least 1000 bootstrap replicates. To determine the homology of  
13  
14 the Dscam-related genes found in metazoans, we estimated the phylogenies of 217 proteins,  
15  
16 including Dscams encompassing the seventh Ig domains to the end (sDscam regions  
17  
18 encompassing the first Ig domains to end) (Fig. 2). This phylogeny was rooted using the  
19  
20  
21 sequence of the *Nematostella vectensis* Dscam (GenBank:ABAV01020293.1)[50]. We used  
22  
23  
24 the duplicated cassettes encoding the Ig1 domain of *sDscam* $\alpha$  to dissect the evolutionary  
25  
26  
27 relationships between the duplicated cassettes with a tree rooted on the *D. melanogaster*  
28  
29  
30 *Dscam1* duplicated exon 9.1, which encodes the Ig7 domain (Fig. S3a). Similarly, we used  
31  
32  
33 the duplicated cassettes, which encode the Ig1+2 domain of *sDscam* $\beta$  and a tree rooted on the  
34  
35  
36 *D. melanogaster Dscam1* exon 9.1, 10, and 11, which encode the Ig7+8 domains (Fig. S3b).  
37  
38  
39

## 40 **Protein three-dimensional structure modelling**

41  
42  
43  
44 All three-dimensional proteins structures were acquired using the Swiss-Model (automated  
45  
46  
47 mode) ([www.swissmodel.expasy.org](http://www.swissmodel.expasy.org)). The structures were displayed and processed using the  
48  
49  
50 PyMOL software package ([www.pymol.com](http://www.pymol.com)).  
51  
52  
53

## 54 **Analysis of differential and biased expression**

55  
56  
57  
58 The RNA-seq data from publically available samples were used to analyse the expression of  
59  
60  
61

nonclassical *Dscam* genes at various developmental stages, tissues types, and cell lines (Additional file 4: Table S3). For each sample, we calculated the RPKM value of the constant exonic region to measure the expression level of each *Dscam* gene from the replicates. The alternative exon encoding Ig7 was selected to calculate the expression level from the replicates for each 5' variable cassette. Considering the short length of the alternative exons, the RNA-seq reads were split into 25-nucleotide (nt) fragments for mapping by using Bowtie 2 software [45], and only the perfectly mapped fragments were retained for expression level calculation. Furthermore, the read counts of a 25-nt fragment with multiple loci were divided by the number of loci, and then assigned equally to each locus for expression level calculation. To eliminate influences on calculations of the expression levels from identical sequences among exon duplicates, both the 25-nt fragments and the full-length RNA-seq reads (150 nt) were used to calculate the expression profiles as previously described [25]. If, as a result of originating from a repetitive region, one RNA-seq read or fragment mapped onto several loci, we divided the RPM value of this RNA-seq read by the number of repetitive loci, then evenly assigned it to each transcript in the expression level calculation.

## Statistical analysis

We used an independent sample *t*-test to assess the relationship between the 5' cassette duplicate and the sDscam gene duplicate. We compared the numbers of 5' cassettes and gene duplicates between the sDscam $\alpha$  and sDscam $\beta$  groups using a two-tailed Student's *t*-test. We defined correlation and significance levels for the number of 5' cassette duplicates and the repeat size in terms of a simple regression model. We also used a two-tailed Student's *t*-test to



compare the differences between the expression levels in groups of mDscams and sDscams across various species. Effects were considered to be statistically significant when  $p < 0.05$ .

## Additional files

Additional file 1: Table S1. Chelicerata species and their genome sources in this study

Additional file 2: Table S2. A list of *Dscam* homologues identified in Chelicerata species.

Additional file 3: Supplementary Figures S1–7.

Additional file 4: Table S3. RNA-seq datasets in Chelicerata species

## Abbreviations

BLAST: Basic local alignment search tool

cDNA: Complementary DNA

Dscam: Down syndrome cell adhesion molecule

EC: Extracellular cadherin domains

EGFP: Enhanced green fluorescent protein

FNIII: Fibronectin III

Ig: Immunoglobulin

LDscam: Large Dscam

1 MDscam: Middle Dscam  
2  
3

4 MEGA: Molecular evolutionary genetics analysis  
5  
6  
7

8 mRNA: Messenger RNA  
9  
10

11  
12 Pcdhs: Protocadherins  
13  
14  
15

16 Pdb ID: Protein data bank identification  
17  
18  
19

20 RNA-seq: RNA Sequencing  
21  
22  
23

24 RPKM: Reads per kilobase of transcript per million mapped reads  
25  
26  
27

28 RPM: Reads per million  
29  
30  
31

32 sDscam: Shortened Dscam  
33  
34  
35

36 3D: Three-dimensional  
37  
38  
39

40 TM: Transmembrane  
41  
42  
43

44 *Mesobuthus martensii*: Mma  
45  
46  
47

48 *Ixodes scapularis*: Isc  
49  
50  
51

52 *Stegodyphus mimosarum*: Smi  
53  
54  
55

56 *Parasteatoda tepidariorum*: Pte  
57  
58  
59  
60  
61  
62  
63  
64  
65

*Metaseiulus occidentalis*: Moc

*Limulus polyphemus*: Lpo

## **Declarations**

## **Ethics approval and consent to participate**

Not applicable.

## **Consent for publication**

Not applicable.

## **Fundings**

This work was supported by research grants from the National Natural Science Foundation of China (31630089, 31430050, 91740104).

## **Author contributions**

YJ conceived of this project. GC, YS, HL, JZ, HM, HD, SC and YW annotated and analyzed the sequences; YJ, GC and SH designed the experiments; SH, GC, PG and WH cloned the nucleotide sequences; YS and JZ conducted phylogenetic and evolutionary analysis; SH and GC conducted alternative splicing analyses; SH and YM analyzed exon junction; YM and GC performed performed the computational analysis of expression level; GC and JZ analyzed promoter activity. YJ and YM performed data correlations. YJ, SH, GC, YM, YS, BX and FS

analyzed the data; YJ and YM wrote the manuscript; all authors discussed the results and commented on the manuscript.

## Availability of data and materials

All data generated or analysed during this study are included in this published article and its Additional files. The *Dscam* gene sequences were deposited into GenBank with accession numbers: MF106020–MF106053; MF066898–MF066905; KX555546–KX555560; KU641421–KU641424.

## Acknowledgments

We thank Dr. Qi zhou, Institute of Life Sciences, Zhejiang University, China, for the critical reading of the manuscript.

## Competing interests

The authors declare that they have no competing interests.

## References

1. Nilsen TW, Graveley BR. Expansion of the eukaryotic proteome by alternative splicing. *Nature*. 2010;463:457-63.
2. Keren H, Lev-Maor G, Ast, G. Alternative splicing and evolution: diversification, exon definition and function. *Nat Rev Genet*. 2010;11:345-55.
3. Schmucker D, Clemens JC, Shu H, Worby CA, Xiao J, Muda M, et al. *Drosophila* Dscam is an axon guidance receptor exhibiting extraordinary molecular diversity. *Cell*.

- 2000;101:671-84.
4. Zhan XL, Clemens JC, Neves G, Hattori D, Flanagan JJ, Hummel T, et al. Analysis of Dscam diversity in regulating axon guidance in *Drosophila* mushroom bodies. *Neuron*. 2004;43:673-86.
  5. Watson FL, Püttmann-Holgado R, Thomas F, Lamar DL, Hughes M, Kondo M, et al. Extensive diversity of Ig-superfamily proteins in the immune system of insects. *Science*. 2005;309:1874-8.
  6. Chen BE, Kondo M, Garnier A, Watson FL, Püttmann-Holgado R, Lamar DR, et al. The molecular diversity of Dscam is functionally required for neuronal wiring specificity in *Drosophila*. *Cell*. 2006;125:607-20.
  7. Dong YM, Taylor HE, Dimopoulos G. AgDscam, a hypervariable immunoglobulin domain-containing receptor of the *Anopheles gambiae* innate immune system. *PLoS Biol*. 2006;4:1137-46.
  8. Soba P, Zhu S, Emoto K, Younger S, Yang SJ, Yu HH, et al. *Drosophila* sensory neurons require Dscam for dendritic self-avoidance and proper dendritic field organization. *Neuron*. 2007;54:403-16.
  9. Wojtowicz WM, Wu W, Andre I, Qian B, Baker D & Zipursky SL. A vast repertoire of Dscam binding specificities arises from modular interactions of variable Ig domains. *Cell*. 2007;130:1134-45.
  10. Hattori D, Chen Y, Matthews BJ, Salwinski L, Sabatti C, Grueber WB, et al. Robust discrimination between self and non-self neurites requires thousands of Dscam1 isoforms. *Nature*. 2009;461:644-648.

11. Miura SK, Martins A, Zhang KX, Graveley BR, Zipursky SL. Probabilistic splicing of Dscam1 establishes identity at the level of single neurons. *Cell*. 2013;155:1166-77.
12. Sun W, You XT, Gogol-Döring A, He HH, Kise Y, Sohn M, et al. Ultra-deep profiling of alternatively spliced *Drosophila* Dscam isoforms by circularization-assisted multi-segment sequencing. *EMBO J*. 2013;32:2029-38.
13. He HH, Kise Y, Izadifar A, Urwyler O, Ayaz D, Parthasarathy A, et al. Cell-intrinsic requirement of Dscam1 isoform diversity for axon collateral formation. *Science*. 2014;344:1182-6.
14. Schmucker D, Chen B. Dscam and DSCAM: complex genes in simple animals, complex animals yet simple genes. *Genes Dev*. 2009;23:147-56.
15. Zipursky SL, Sanes JR. Chemoaffinity revisited: dscams, protocadherins, and neural circuit assembly. *Cell*. 2010;143:343-53.
16. Lefebvre JL, Kostadinov D, Chen WV, Maniatis T, Sanes JR. Protocadherins mediate dendritic self-avoidance in the mammalian nervous system. *Nature*. 2012;488:517-21.
17. Zipursky SL, Grueber WB. The molecular basis of self-avoidance. *Annu. Rev. Neurosci*. 2013;36:547-68.
18. Wu Q, Maniatis T. A striking organization of a large family of human neural cadherin-like cell adhesion genes. *Cell*. 1999;97:779-90.
19. Tasic B, Nabholz CE, Baldwin KK, Kim Y, Rueckert EH, Ribich SA, et al. Promoter choice determines splice site selection in protocadherin alpha and -gamma pre-mRNA splicing. *Mol. Cell*. 2002;10:21-33.
20. Wang XZ, Su H, Bradley A. Molecular mechanisms governing Pcdh-gamma gene

- expression: Evidence for a multiple promoter and cis-alternative splicing model. *Genes Dev.* 2002;16:1890-905.
21. Chen WV, Nwakeze CL, Denny CA, O'Keeffe S, Rieger MA, Mountoufaris G, et al. *Pcdhac2* is required for axonal tiling and assembly of serotonergic circuitries in mice. *Science.* 2017;356:406-11.
22. Mountoufaris G, Chen WV, Hirabayashi Y, O'Keeffe S, Chevee M, Nwakeze CL, et al. Multicluster *Pcdh* diversity is required for mouse olfactory neural circuit assembly. *Science.* 2017;356:411-4.
23. Armitage SA, Freiburg RY, Kurtz J, Bravo IG. The evolution of *Dscam* genes across the arthropods. *BMC Evol Biol.* 2012;12:53-67.
24. Brites D, Brena C, Ebert D, Du Pasquier L. More than one way to produce protein diversity: duplication and limited alternative splicing of an adhesion molecule gene in basal arthropods. *Evolution.* 2013;67:2999-3011.
25. Yue Y, Meng YJ, Ma HR, Hou SQ, Cao GZ, Hong WL, et al. A large family of *Dscam* genes with tandemly arrayed 5' cassettes in Chelicerata. *Nat Commun.* 2016;7:11252.
26. Lee MS, Soubrier J, Edgecombe GD. Rates of phenotypic and genomic evolution during the Cambrian explosion. *Curr. Biol.* 2013;23:1889-95.
27. Leite DJ, McGregor AP. Arthropod evolution and development: recent insights from chelicerates and myriapods. *Curr Opin Genet Dev.* 2016;39:93-100.
28. Kopelman NM, Lancet D, Yanai I. Alternative splicing and gene duplication are inversely correlated evolutionary mechanisms. *Nat Genet.* 2005;37:588-9.
29. Gulia-Nuss M, Nuss AB, Meyer JM, Sonenshine DE, Roe RM, Waterhouse RM, et al.

- 1 Genomic insights into the *Ixodes scapularis* tick vector of Lyme disease. Nat Commun.  
2  
3 2016;7:10507.  
4  
5  
6 30. Yamakawa K, Huot YK, Haendelt MA, Hubert R, Chen XN, Lyons GE, et al. DSCAM:  
7  
8 a novel member of the immunoglobulin superfamily maps in a Down syndrome region  
9  
10 and is involved in the development of the nervous system. Hum. Mol. Genet.  
11  
12 1998;7:227-37.  
13  
14  
15  
16  
17 31. Celotto AM, Graveley BR. Alternative splicing of the Drosophila Dscam pre-mRNA is  
18  
19 both temporally and spatially regulated. Genetics. 2001;159:599-608.  
20  
21  
22  
23 32. Neves G, Zucker J, Daly M, Chess A. Stochastic yet biased expression of multiple  
24  
25 Dscam splice variants by individual cells. Nat. Genet. 2004;36:240-6.  
26  
27  
28 33. Schreiner D, Weiner JA. Combinatorial homophilic interaction between  
29  
30 gamma-protocadherin multimers greatly expands the molecular diversity of cell adhesion.  
31  
32 Proc Natl Acad Sci U S A. 2010;107:14893-8.  
33  
34  
35  
36 34. Goodman KM, Rubinstein R, Thu CA, Bahna F, Mannepalli S, Ahlsén G, et al. Structural  
37  
38 basis of diverse homophilic recognition by clustered  $\alpha$ - and  $\beta$ -protocadherins. Neuron.  
39  
40 2016;90:709-23.  
41  
42  
43  
44 35. Shapiro L, Love J, Colman DR. Adhesion molecules in the nervous system: structural  
45  
46 insights into function and diversity. Annu Rev Neurosci. 2007;30:451-74.  
47  
48  
49  
50 36. Ravi V, Yu WP, Pilla NE, Lian MM, Tay BH, Tohari S, et al. Cyclostomes lack clustered  
51  
52 protocadherins. Mol Biol Evol. 2016;33:311-5.  
53  
54  
55  
56 37. Hulpiau P, van Roy F. New insights into the evolution of metazoan cadherins. Mol Biol  
57  
58 Evol. 2011;28:647-57.  
59  
60  
61  
62  
63  
64  
65



- 1 38. Guo Y, Xu Q, Canzio D, Shou J, Li JH, Gorkin DU, et al. CRISPR inversion of CTCF  
2  
3 sites alters genome topology and enhancer/promoter function. *Cell*. 2015;162:900-10.  
4  
5  
6 39. Hozumi N, Tonegawa S. Evidence for somatic rearrangement of immunoglobulin genes  
7  
8 coding for variable and constant regions. *Proc Natl Acad Sci*.1976;73:3628-32.  
9  
10  
11 40. Tonegawa S. Somatic generation of antibody diversity. *Nature*. 1983;302:575-81.  
12  
13  
14 41. Du Pasquier L, Zucchetti I, Santis, RD. Immunoglobulin superfamily receptors in  
15  
16 protochordates: before RAG time. *Immunol. Rev*.2004;198:233-48.  
17  
18  
19 42. Kurtz J, Armitage SAO. Alternative adaptive immunity in invertebrates. *Trends in*  
20  
21 *Immunology*. 2006;27:493-6.  
22  
23  
24 43. Sanggaard KW, Bechsgaard JS, Fang X, Duan J, Dyrland TF, Gupta V, et al. Spider  
25  
26 genomes provide insight into composition and evolution of venom and silk. *Nat*  
27  
28 *Commun*. 2014;5:3765.  
29  
30  
31 44. Cao ZJ, Yu Y, Wu YL, Hao P, Di ZY, He YW, et al. The genome of *Mesobuthus martensii*  
32  
33 reveals a unique adaptation model of arthropods. *Nat Commun*. 2013;4:2602.  
34  
35  
36 45. Langmead B, Salzberg SL. Fast gapped-read alignment with Bowtie 2. *Nat methods*.  
37  
38 2012;9:357-9.  
39  
40  
41 46. Mitchell A, Chang HY, Daugherty L, Fraser M, Hunter S, Lopez R, et al. The InterPro  
42  
43 protein families database: the classification resource after 15 years. *Nucleic Acids Res*.  
44  
45 2015;43:D213-21.  
46  
47  
48 47. Sigrist CJ, de Castro E, Cerutti L, Cucho BA, Hulo N, Bridge A, et al. New and  
49  
50 continuing developments at PROSITE. *Nucleic Acids Res*. 2013;41 (Database  
51  
52 issue):D344-7.  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

48. Larkin MA, Blackshields G, Brown NP, Chenna R, McGettigan PA, McWilliam H, et al. Clustal W and Clustal X version 2.0. *Bioinformatics*. 2007;23:2947-8.
49. Kumar S, Stecher G, and Tamura K. MEGA7: Molecular evolutionary genetics analysis version 7.0 for bigger datasets. *Mol Biol Evol*. 2016; 33:1870-4.
50. Ganot P, Zoccola D, Tambutté E, Voolstra CR, Aranda M, Allemand D, et al. Structural molecular components of septate junctions in cnidarians point to the origin of epithelial junctions in eukaryotes. *Mol Biol Evol*. 2015; 32:44-62.
51. Regier JC, Shultz JW, Zwick A, Hussey A, Ball B, Wetzer R, et al. Arthropod relationships revealed by phylogenomic analysis of nuclear protein-coding sequences. *Nature*. 2010;463:1079-83.

## Figure legends

### Fig. 1. Genome-wide identification and classification of *Dscams* in Chelicerata. (a)

Schematic representation of Dscam structures in Chelicerata. Ig, immunoglobulin domains; FNIII, fibronectin III domains. The N-terminal small boxes represent the leader peptides. The black and green boxes represent the transmembrane (TM) and cytoplasmic domains. Five *Dscam* types (*LDscam*, *mDscam*, *sDscam $\alpha$* , *sDscam $\beta$* , and *sDscam $\gamma$* ) are classified based on their size and structure. *LDscam* shares structures identical to classical *Dscam*. *mDscam* lacks Ig10 and FNIII 3–4 domains of classical *Dscam*. *sDscam $\alpha$*  contains variable N-terminal IgI (blue), which corresponds to the variable Ig7 domain of *Drosophila Dscam1*. *sDscam $\beta$*  contains variable N-terminal Ig1+2 domains (coloured), which correspond to the variable Ig7+8 domains of *Drosophila Dscam1*. Numbers in parentheses refer to the numbers of 5'

variable cassettes in *Metaseiulus occidentalis* *sDscama* and *sDscamβ2*. *sDscamγ* shares domains similar to *sDscama* and *sDscamβ*, albeit with no tandemly arrayed cassettes in the 5' regions. Numbers in parentheses refer to the numbers of tandem cassettes. (b) Phylogenetic distribution of *Dscams* and isoform members in chelicerates. *Dscams* are shown associated with a cladogram of phylogenetic relationships in this study [51]. # Indicates the putative numbers of *Dscam* isoforms in various arthropod species caused by either gene and/or exon duplication, estimated by the number of Ig7 or orthologues.

**Fig. 2. The evolutionary relationships and protein structures of arthropoda *Dscams*.** The tree is based on amino acid sequence alignment of *Dscams* encompassing the seventh Ig domains to end (*sDscam* regions encompassing the first Ig domains to end), and is rooted using the sequence of the *Nematostella vectensis* *Dscam* (GenBank:ABAV01020293.1)[50]. The support values at the nodes are bootstrap values relative to 1000 replicates. We considered five *Dscam* types (*LDscam*, *mDscam*, *sDscama*, *sDscamβ*, and *sDscamγ*) in Chelicerates, which are detailed in Table S2 (Additional file 2). *L. vannamei* (Lva) *Dscam1* (GQ154653), *D. melanogaster* (Dme) *Dscam1–4* (CG17800; CG42256; CG31190; CG42330), *A. mellifera* (Ame) *Dscam* (AAT96374; BAF03050.1; XM\_396307), *D. pulex* *Dscam1* (EU307884), *A. californica* (Aca) *Dscam* (ABS30432.1), *S. kowalevskii* (Sko) *Dscam* (XP\_006825869.1), *S. purpuratus* (Spu) *Dscam* (XP\_011665747.1), *D. japonica* (Dja) *Dscam* (BAE94189.1), and *H. sapiens* (Hsa) *Dscam-L* (AAL57166.1) are included. For other canonical *Dscam* sequences from *S. maritima*, refer to recent references [24]. We collapsed the monophyletic clades of *sDscams* and *mDscams* for visualization convenience. Their detailed phylogenetic relationships are shown in Figures S1 and S2 (Additional file 3),

respectively. The sequential shortening of the Ig and FNIII domains of canonical *Dscam* is indicated by orange squares.

**Fig. 3. Similarities among 5' variable cassettes of *sDscam* genes in representative species from each of the major Chelicerata clades.** (a) Heatmap of pairwise sequence identities of the duplicated cassettes encoding the Ig1 domain of *sDscam* $\alpha$  isoforms from *Metaseiulus occidentalis* (Moc), *Ixodes scapularis* (Isc), *Mesobuthus martensii* (Mma), *Stegodyphus mimosarum* (Smi), and *Limulus polyphemus* (Lpo). To simplify visualisation, the cassette repeats are depicted in the heatmaps in the same linear order in which they reside in each gene. (b) Heatmaps of the percent identities of the duplicated cassettes encoding Ig1+2 domains of *sDscam* $\beta$  in Chelicerata species. The order of the cassettes in the heatmap corresponds to the linear order of each cassette in the genome. (c) Phylogenetic tree for the 268 duplicated cassettes encoding Ig1 domains of *sDscam* $\alpha$ . The support values at the nodes are bootstrap values relative to 1000 replicates. Numbers in parentheses refer to the number of 5' variable cassettes of the corresponding gene. Clades were collapsed for visualisation convenience. (d) Phylogenetic tree for the 243 duplicated cassettes encoding Ig1+2 domains of *sDscam* $\beta$ . Collapsed clades plus crosses are depicted as a combined clade composed of these duplicated cassettes from different species.

**Fig. 4. Inverse relationship between the number of 5' cassette tandem duplicates and gene duplicates.** (a, b) The number of 5' cassette duplicates correlated inversely with the number of gene duplicates in *sDscam* $\alpha$  (a) and *sDscam* $\beta$  (b) in Chelicerata species. (c, d) Comparison of the number of 5' cassette (c) and gene duplicates (d) between *sDscam* $\alpha$  and

*sDscamβ*. Data are expressed as mean  $\pm$  standard deviation (SD). (e–i) The number of 5' cassette tandem duplicates correlated inversely with the duplicate size. This reverse relationship is conserved in *M. occidentalis* (e), *I. scapularis* (f), *M. martensii* (g), *S. mimosarum* (h), and *L. polyphemus* (i). (j) Comparison of the size of cassette repeats in *sDscama* and *sDscamβ*. (k, l) The size of the cassette repeats correlated with the genome size of Chelicerata in *sDscama* (k) and *sDscamβ* (l). The number of 5' cassette tandem duplicates is indicated by the circle size. (m) Intron loss occurred independently in the *sDscamβ1* and *sDscamβ2* genes of *M. occidentalis*. The introns are represented with lines (not drawn to scale). Red circles indicate the intron loss, and the arrows indicate the transcription start sites.

**Fig. 5. Differential expression of nonclassical *Dscam* genes.** (a) Heatmap of expression of 31 *Dscam* genes in various tissues of *M. martensii*. The expression level for each transcript is shown as reads per kilobase of transcript per million mapped reads (RPKM) of its corresponding constitutive exons. The 25-nt fragmented RNA-sequencing datasets were mapped to calculate the relative expression level. The maximum expression levels of *mDscam8* were found in hemocytes, whereas *mDscam10* and *mDscam15* were highly expressed in poison glands. (b) The differing expression patterns of *Dscams* in *I. scapularis* lineages. This result indicates that many *Dscams* were preferentially and specifically expressed in certain lines. Data are expressed as the mean  $\pm$  SD from three independent experiments.

**Fig. 6. *sDscam* genes are biased to be more highly expressed.** (a) In *M. martensii*, *sDscam* genes are biased to be more highly expressed than *mDscam* genes. We calculated p values

using a two-tailed Student's *t*-test. As *LDscams* occurs much less frequently than *mDscams* and *sDscams*, we did not analyze the statistical differences between them. **(b)** *sDscam* genes are biased to be more highly expressed than *mDscam* genes in *I. scapularis*. **(c)** The expression bias is conserved in *S. mimosarum* (Smi, SRR1015314), *P. tepidariorum* (Pte, SRR1824487), *M. occidentalis* (Moc, SRR446504), and *L. polyphemus* (Lpo, SRX1323743).

**Fig. 7. Comparisons of organisation and origin of clustered *Dscam* and *Pcdh* genes. (a)**

Schematic diagram for the *sDscama* gene in Chelicerata. Symbols used are the same as in Figure 1. Each variable cassette was transcribed by an alternative promoter followed by alternative splicing. The variable cassette encoded the N-terminal Ig1 domain (blue). Tertiary structure model of Ig1 of *I. scapularis* *sDscam* is shown on the left. **(b)** Schematic diagram for the *Pcdha* genes in vertebrates. The *Pcdh* gene cluster contains exons that encode 14 extracellular and TM domains [18]. Each repeat is preceded by a promoter, and encodes extracellular and TM domains. Tertiary structure model of the EC1 domain of *Pcdha* is shown on the left, which is similar to that of the Ig1 domain of *sDscam*. **(c, d)** Comparisons between the evolutionary origin and expansion of clustered *Dscam* (c) and *Pcdh* (d) genes. The 5' clustered organisation of both *sDscam* and *Pcdh* genes may originate from the shortening and expansion of the ectodomains of canonical *Dscam* via sequential duplication and mutation.

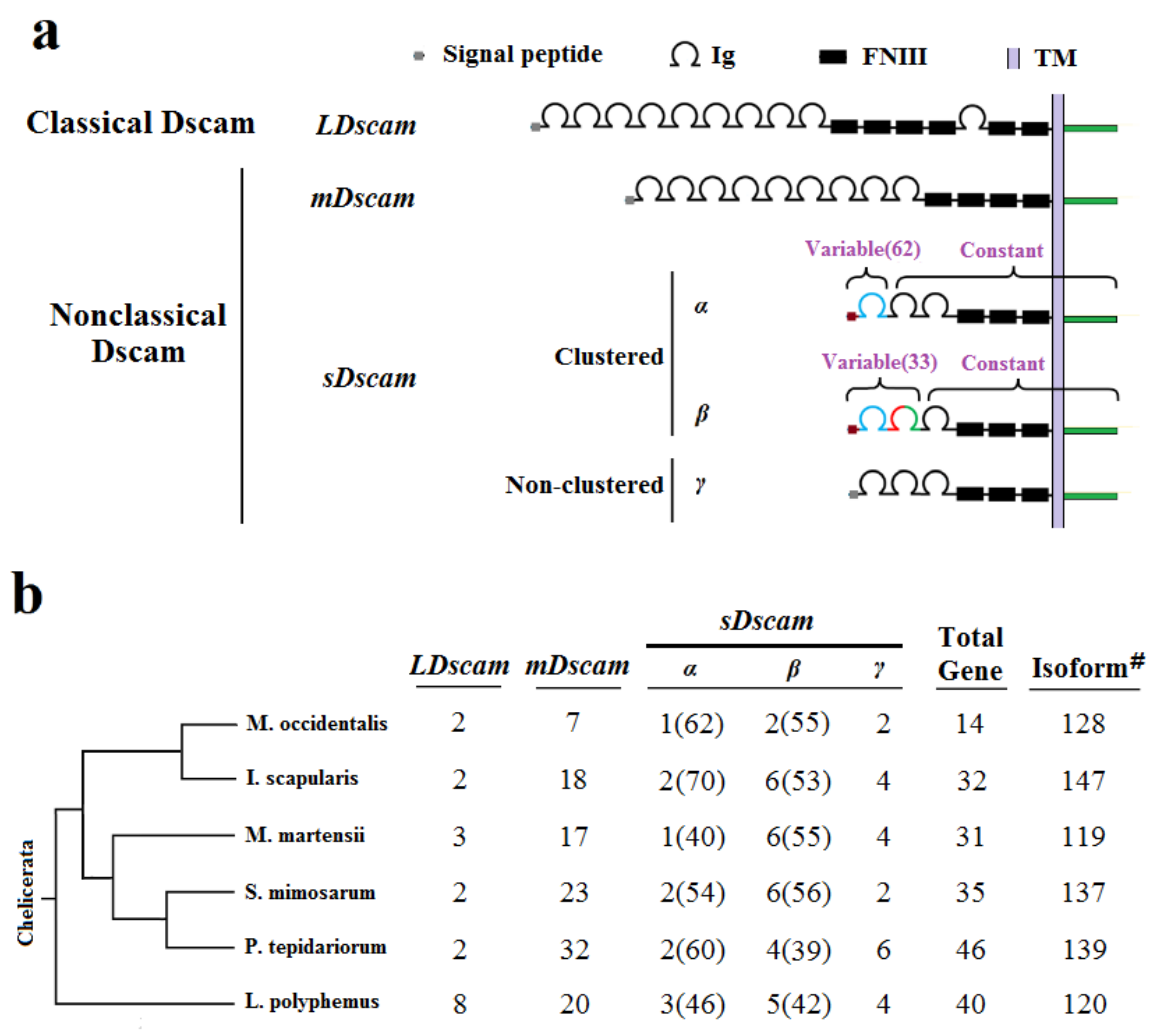


Fig. 1

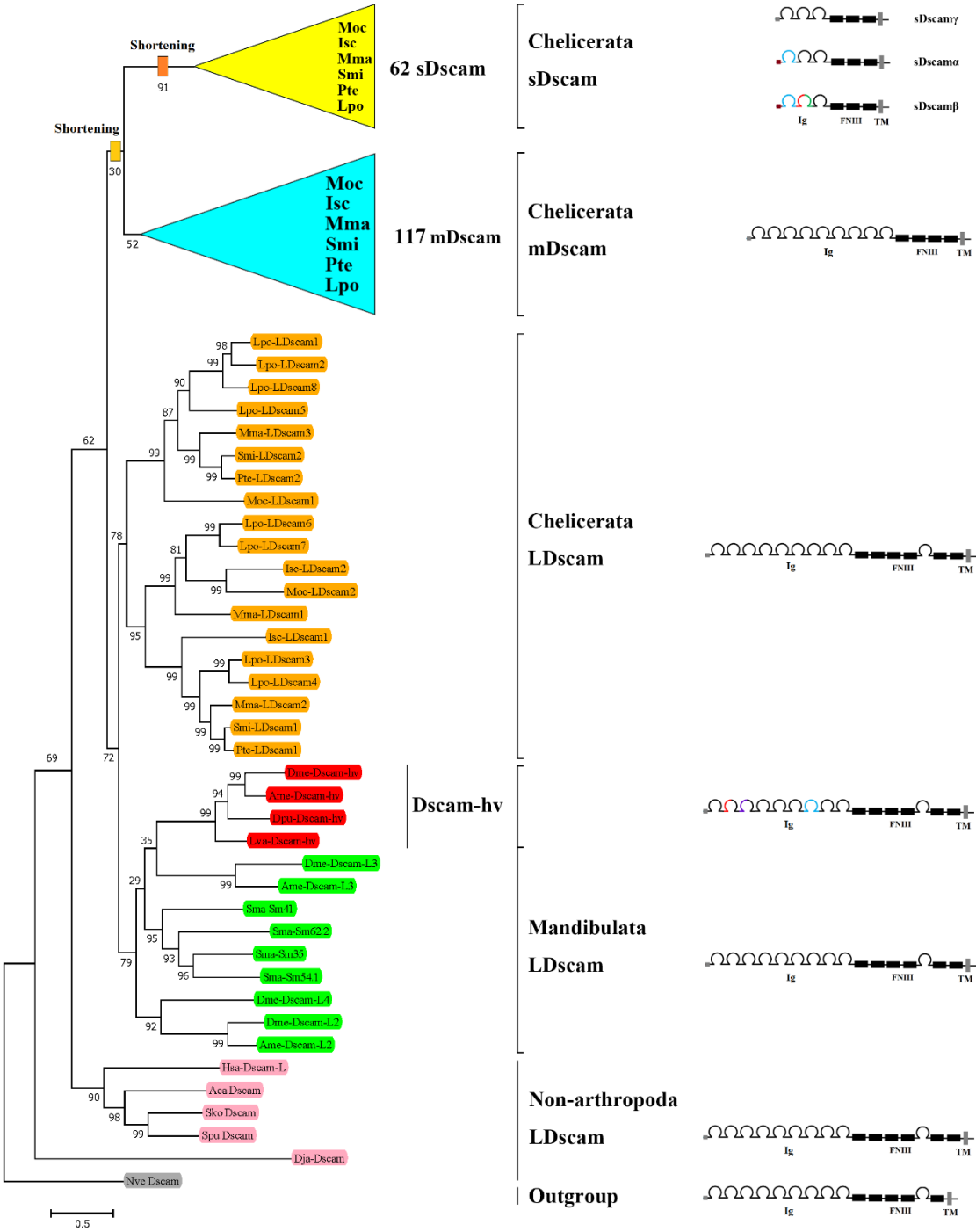


Fig. 2



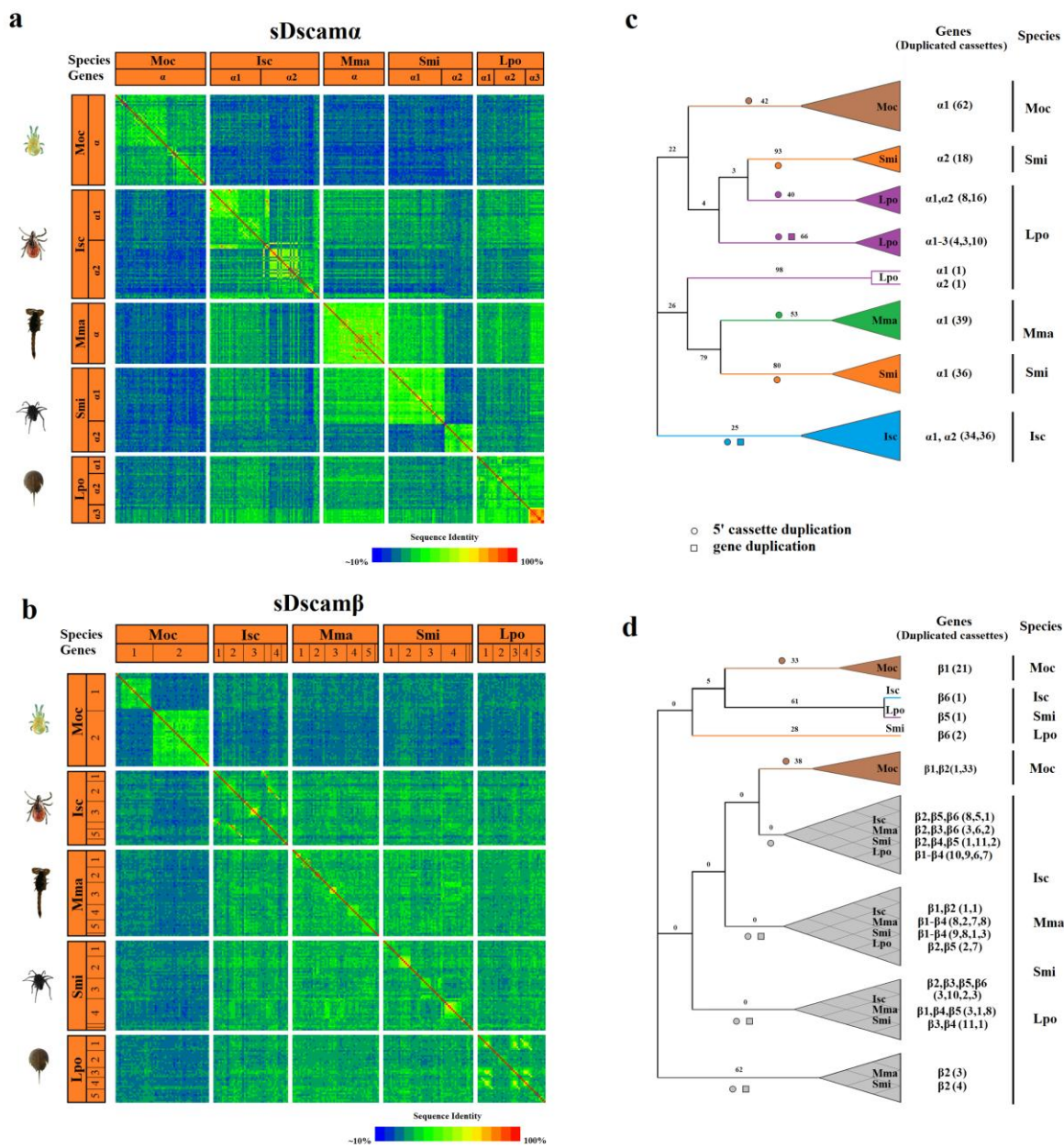


Fig. 3

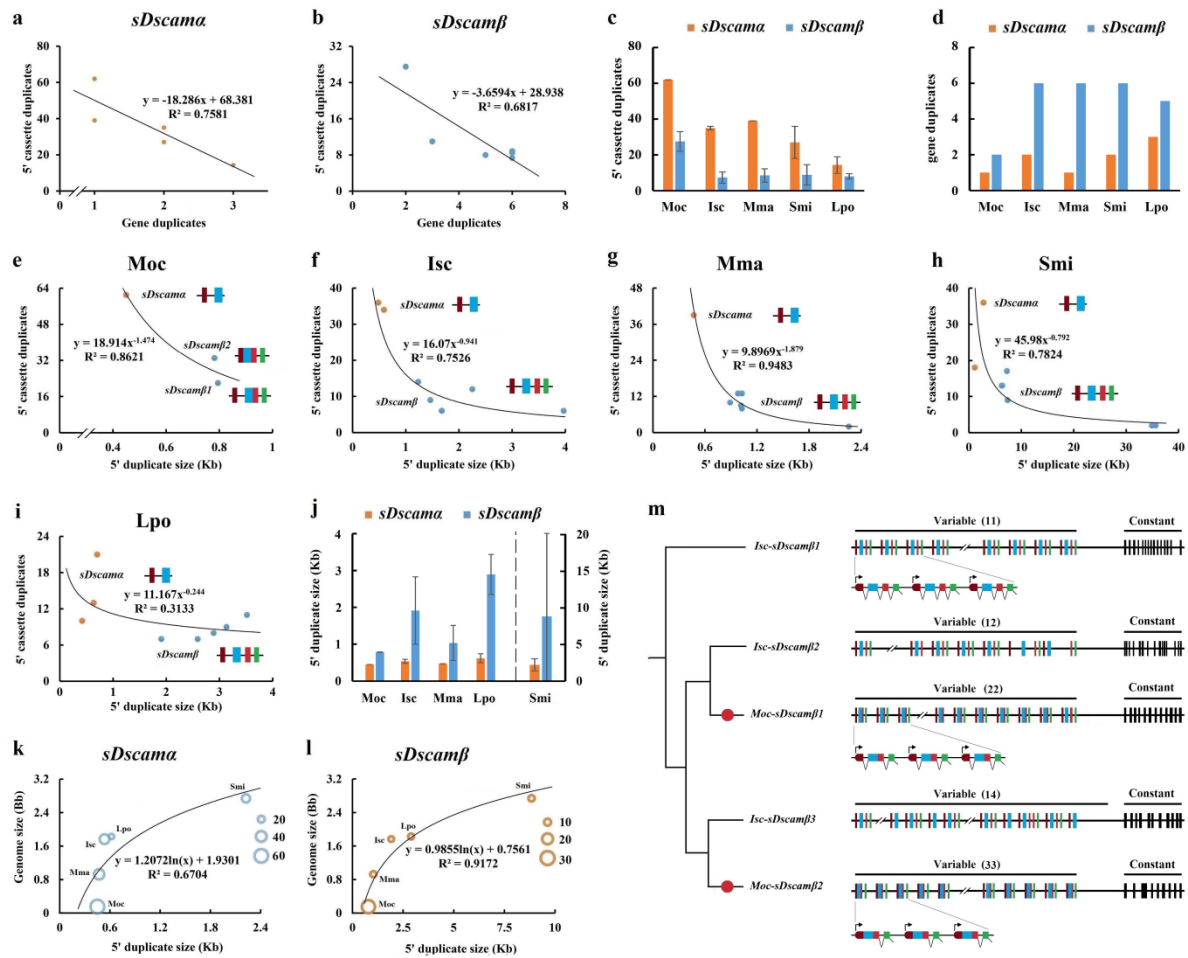


Fig. 4

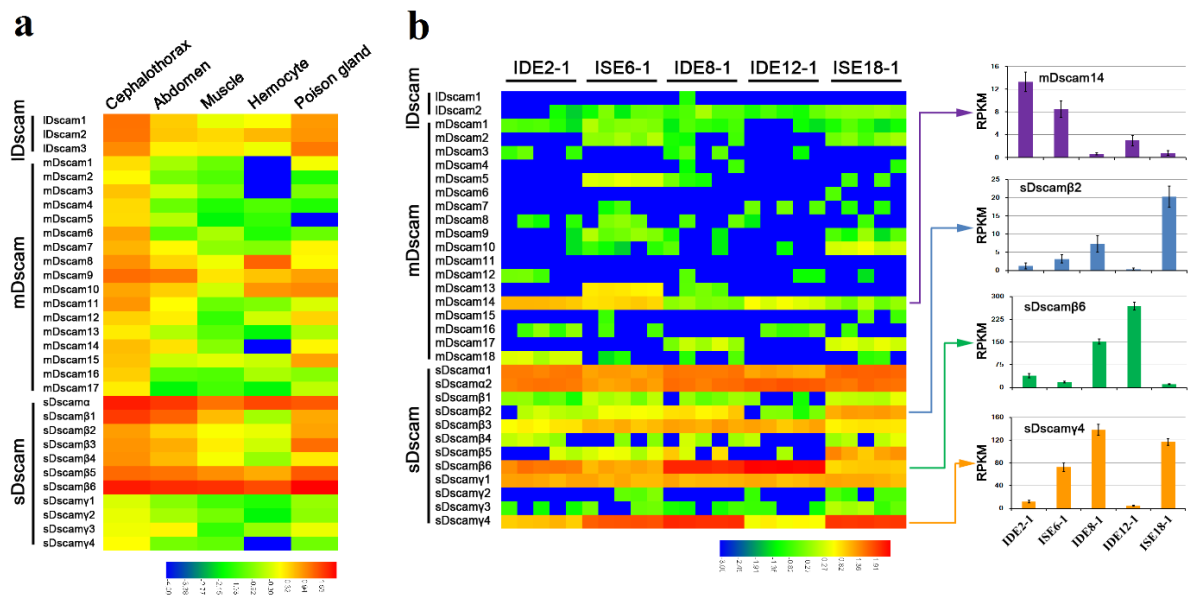


Fig. 5

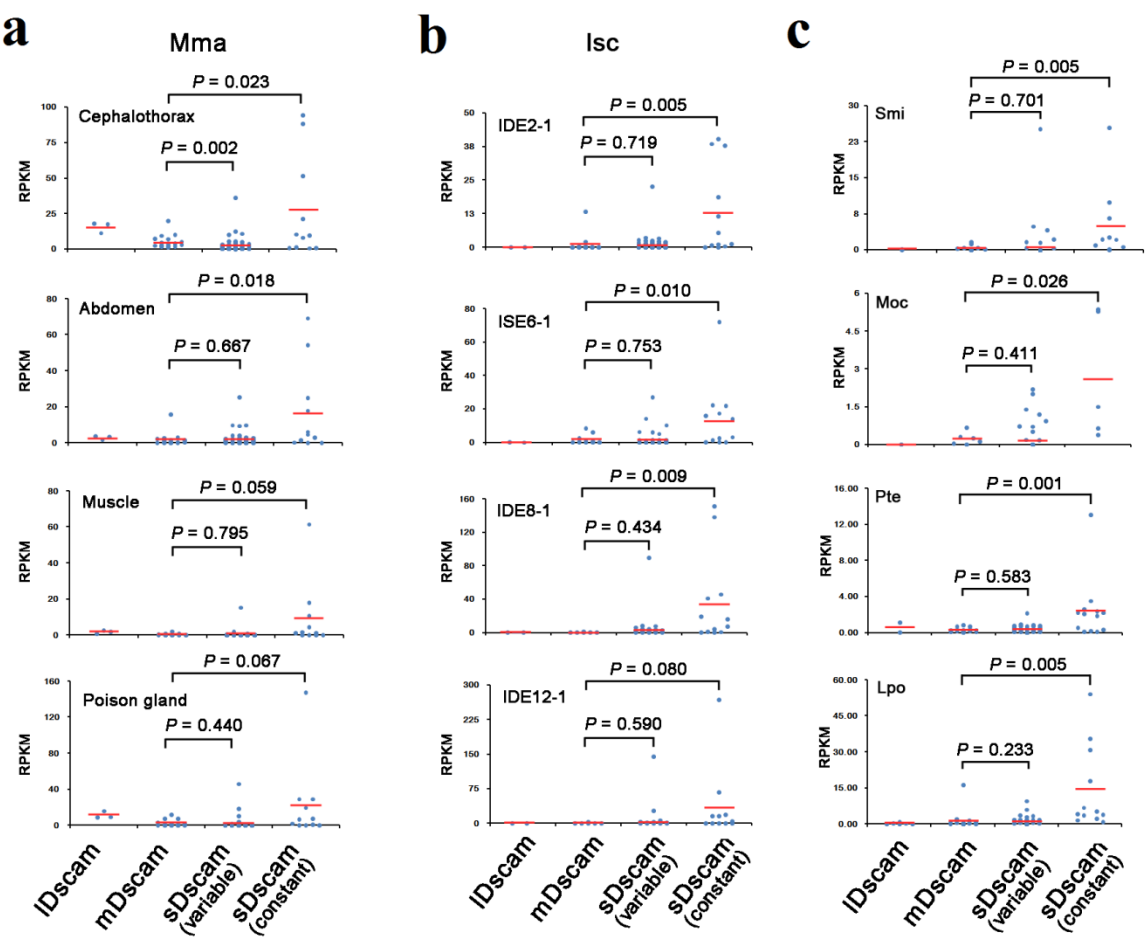
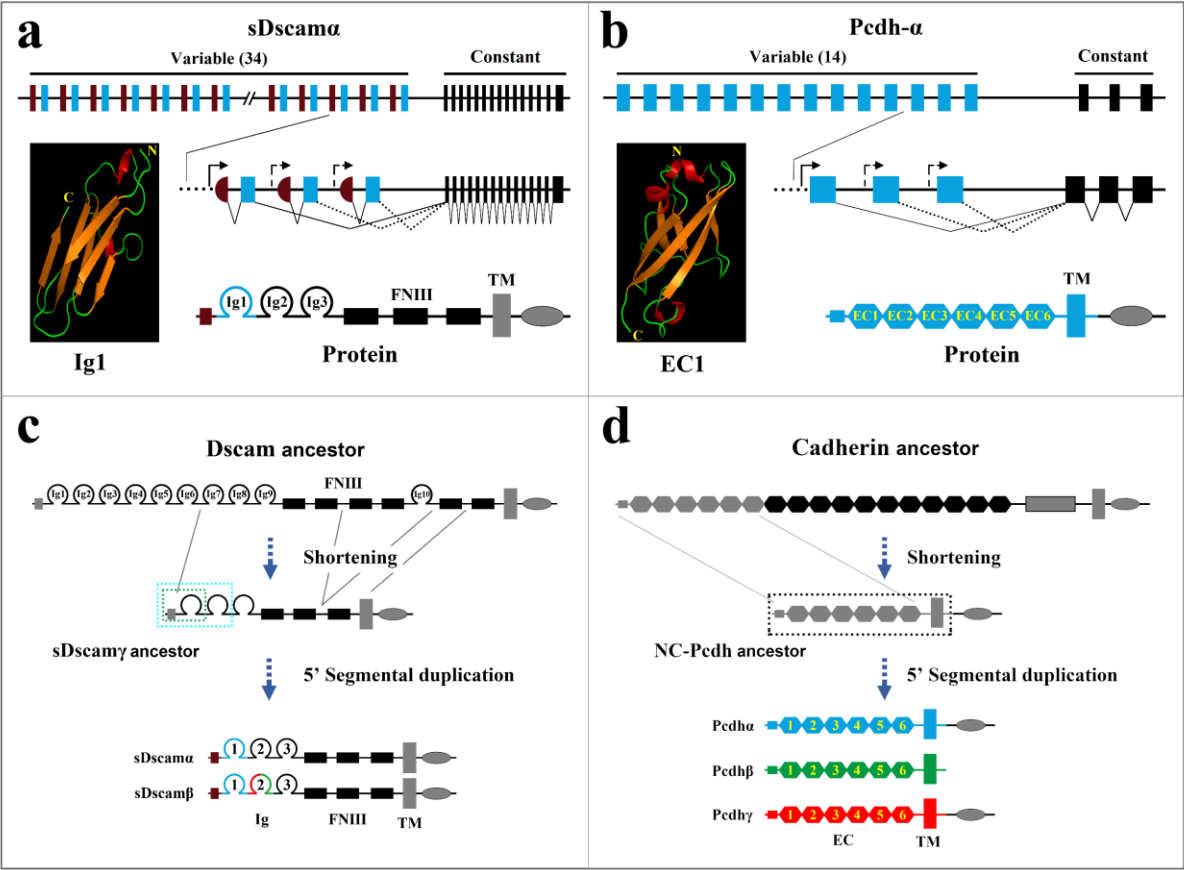


Fig. 6



**Fig. 7**

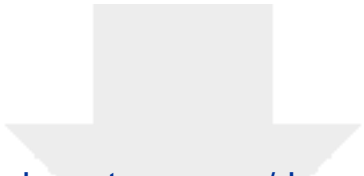


[Click here to access/download](#)

**Supplementary Material**

[Additional file 1-Table S1 Chelicerata species.xlsx](#)

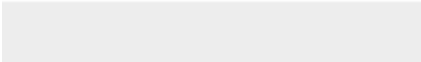





[Click here to access/download](#)

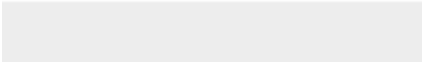

**Supplementary Material**

Revised Additional file 2-Table S2.xlsx





Click here to access/download  
**Supplementary Material**  
renamed\_c95b5.pdf



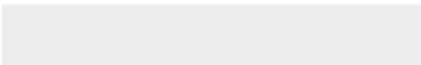
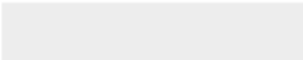


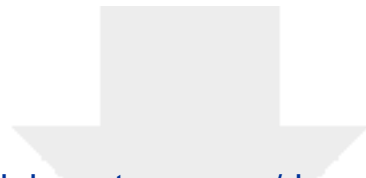


[Click here to access/download](#)

**Supplementary Material**

Additional file 4-Table S3 RNA-seq datasets.xlsx





[Click here to access/download](#)

**Supplementary Material**

JinResponse To Reviewer Comment.pdf

