

Article; MBE Discoveries section

Assortative Mating Drives Linkage-Disequilibrium between Sperm and Egg Recognition  
Protein loci in the Sea Urchin *Strongylocentrotus purpuratus*

Andres Plata Stapper<sup>1,2</sup>, Peter Beerli<sup>3</sup> and Don R Levitan<sup>1,4</sup>

1 Department of Biological Science, Florida State University, Tallahassee, Florida,  
United States of America

2. Current address, Department of Otolaryngology – HNS, Stanford University, Stanford,  
California, United States of America

3. Department of Scientific computing, Florida State University, Tallahassee, Florida,  
United States of America

4. Corresponding Author

Running Head: Assortative mating drives linkage of reproductive loci

## Abstract

Sperm and eggs have interacting proteins on their surfaces that influence their compatibility during fertilization. These proteins are often polymorphic within species, producing variation in gamete affinities. We first demonstrate the fitness consequences of various sperm binding protein (*Bindin*) variants in the sea urchin *Strongylocentrotus purpuratus*, and assortative mating between males and females based on their sperm *Bindin* genotype. This empirical finding of assortative mating based on sperm *Bindin* genotype could arise by linkage-disequilibrium (LD) between interacting sperm and egg recognition loci. We then examine sequence variation in eight exons of the sea urchin egg receptor for sperm *Bindin* (*EBRI*). We find little evidence of LD among the eight exons of *EBRI*, yet strong evidence for LD between sperm *Bindin* and *EBRI* overall, and varying degrees of LD between sperm *Bindin* among the eight exons. We reject the alternate hypotheses of LD driven by shared evolutionary histories, population structure, or close physical linkage between these interacting loci on the genome. The most parsimonious explanation for this pattern of LD is that it represents selection driven by assortative mating based on interactions among these sperm and egg loci. These findings indicate the importance of ongoing sexual selection in the maintenance of protein polymorphisms and LD, and more generally highlight how LD can be used as an indication of current mate choice, as opposed to historic selection.

## Introduction

Gamete recognition proteins (GRPs) mediate the compatibility between sperm and eggs and can influence patterns of reproductive success within species (Palumbi

1999; Levitan and Farrell 2006; Levitan and Stapper 2010; Levitan 2012) and the degree of reproductive isolation across species (Zigler and Lessios 2003, Zigler et al. 2005). Functional variation in gamete recognition proteins within a population, under some conditions, is predicted to result in assortative mating (Gavrilets and Waxman 2002; Haygood 2004; Tomaiuolo and Levitan 2010) and provides a potential mechanism of sympatric speciation (Gavrilets and Waxman 2002). Assortative mating is here defined as variation in gamete compatibilities that drive non-random fusion of sperm and eggs. Assortative mating based on matching of specific alleles from two loci (one expressed on sperm and the other eggs) produce offspring with non-random associations of alleles across loci. This matching is predicted to result in linkage disequilibrium across sperm and egg recognition loci (Payne and Krakauer 1997; Dobelli 2005; Clark et al. 2009; Tomaiuolo and Levitan 2010).

Patterns of sperm availability are predicted to impose differential selection on gamete recognition proteins that might drive protein polymorphism and the degree of linkage disequilibrium (Tomaiuolo and Levitan 2010). Conditions of sperm limitation are predicted to select for the most compatible sperm and egg recognition proteins and select against alternate variants with lower compatibility. This would lead to purifying selection and reduced variation in these proteins; without variation, linkage disequilibrium would not occur. Alternately, conditions of sperm over-abundance can result in polyspermy and developmental failure. In such cases, protein variants in the egg (and in some conditions the sperm) with lower gametic affinities can be favored as they decrease the rate of fertilization and allow blocks to polyspermy to become effective. These conditions lead to polymorphisms in egg proteins, which allow novel sperm

proteins that match these new egg variants to persist in the population. The more restrictive conditions in which sperm protein variants with lower affinities can be selected involve gamete interactions in which eggs are surrounded by overly abundant sperm from a single male. In this scenario, in which males do not directly compete for single eggs, reduced gamete affinities can be selected in either sperm or egg recognition proteins. These conditions that result in a balanced polymorphism of protein variants in gamete recognition proteins are predicted to result in linkage-disequilibrium between the matched sets of sperm and egg protein variants (Tomaiuolo and Levitan 2010). Prior work has demonstrated conditions of sperm limitation select for sperm protein variants with high affinities whereas conditions of polyspermy select for sperm protein variants with lower affinities (Levitan and Ferrell 2006; Levitan and Stapper 2010), and that shifts in population density and the degree of sperm limitation in natural populations have led to contemporary evolution of these proteins in a direction predicted by genotype-specific gamete affinities (Levitan 2012).

Indirect evidence of linkage disequilibrium driven by assortative mating has come from studies examining the reproductive success of sea urchins as a function of their sperm *Bindin* protein. These studies indicate that not only can sperm *Bindin* genotype predict male reproductive success, but also female reproductive success (Palumbi 1999; Levitan and Ferrell 2006; Levitan and Stapper 2010) and further that males and females that match at sperm *Bindin* appear to have higher compatibility at fertilization than mismatched mates (Palumbi 1999; Levitan and Ferrell 2006). Direct evidence of co-evolution and linkage disequilibrium between sperm and egg loci has been detected in abalone (Clark et al 2009), but without the matching data on the assortative mating that

might drive this relationship. Here we provide evidence of assortative mating in the purple sea urchin *Strongylocentrotus purpuratus*, examine patterns of variation in sperm *Bindin* and the egg receptor for sperm *Bindin* (*EBR1*), and test for linkage disequilibrium between these sperm and egg loci in a novel way using the program Migrate (Beerli 2006; Beerli and Palczewski 2010). We examine and reject the alternative hypotheses that linkage disequilibrium is caused by shared evolutionary histories or the genomic proximity of these interacting loci.

#### *Sperm Bindin and EBR1 proteins*

Gamete recognition proteins often show a molecular signature of very rapid evolution (Metz et al. 1998; Swanson and Vaquier 2002; Clark et al. 2006; Turner and Hoekstra 2008). However, identifying and studying interacting GRPs has been difficult, and only a few potentially informative complexes of interacting GRPs have been identified, such as the *Lysin-Verl* complex in abalone (Swanson and Vacquier 1998; Galindo et al. 2003; Aagaard et al. 2010) and the *Bindin-EBR1* complex in echinoderms (Kamei and Glabe 2003; Pujolar and Pogson 2011; Hart 2013). Although from a molecular perspective the *Bindin-EBR1* system is not as well characterized as the *Lysin-Verl* complex, it is potentially highly informative because sperm *Bindin* is often polymorphic within species, and these protein variants have known fitness consequences based on patterns of sperm availability (Levitan and Ferrell 2006, Levitan and Stapper 2010; Levitan 2012).

The sea urchin sperm *Bindin* protein is expressed in acrosome-reacted sperm, and has been isolated and sequenced at the protein level (Vaquier 1977). The protein is only expressed in sea urchin sperm late in spermatogenesis (Nisioka et al. 1990; Cameron et

al. 1990), and has been shown to be a testis-specific transcript (Gao et al. 1986) not expressed in females (Minor et al. 1989). This protein is highly variable across species and the degree of protein divergence predicts the degree of gametic compatibility (Zigler et al. 2005).

*EBRI* is a large glycoprotein (350kDa) expressed on the surface of sea urchin eggs. It is thought to mediate species-specific egg and sperm fusion via its interaction with sperm *Bindin* (Kamei and Glabe 2003), and consists of a disintegrin and metalloprotease, reprotysin and an ADAM cis-rich domain followed by a sequence of repetitive domains: 17 *Tsp*, 10 *Cub* and 9 *Hyr* (Figure 1). The repetitive part of the protein is organized as 9 *Tsp* repeats followed by eight and a half tandem *Tsp/Cub* repeats (*Tsp* -Thrombospondin type1 repeat-, and *Cub* - a domain identified in complement proteins-) and finally by a species-specific domain consisting of hyalin-like (*Hyr*) repeats; *Hyr* is a motif that has been predicted to be involved in cellular adhesion (Callebaut et al. 2000) and has an immunoglobulin structure. Interestingly, responsibility for species specificity between *Strongylocentroid* species has been attributed to the expression of *Hyr* in *S. purpuratus*, because these repeats were not detected in the transcriptomic analysis of the red urchin *S. franciscanus* (Kamei and Glabe 2003).

Little is known about the intra- and interspecific diversity and selection patterns of *EBRI* and the coevolutionary dynamics of sperm *Bindin* and *EBRI* in sea urchins. A recent study by Pujolar and Pogson (2011), which targeted a small region 186bp long (exon 22), corresponding to the 8<sup>th</sup> *Tsp* repeat in *S. franciscanus*, *S. droebachiensis*, *S. purpuratus* and *S. pallidus*, showed a low mean dN/dS ratio at this small region ( $\omega = 0.26$ ), suggesting purifying selection, but pinpointed a single amino acid with a high

dN/dS ratio ( $\omega = 7.78$ ). The authors also showed that the small region of *EBR1* targeted had elevated rates of replacement substitutions only in *S. purpuratus* and *S. pallidus*, while no evidence of positive selection was identified for *S. franciscanus* and *S. droebachiensis*, suggesting that *EBR1* is under different selective pressures in different stronglycentroid species (Pujolar and Pogson 2011). In their study, Pujolar and Pogson were unable to establish the coevolutionary dynamics between *Bindin* and *EBR1*, despite finding an expected signal of higher positive selection in *Bindin* than in *EBR1*; this could occur if the small *EBR1* target region sequenced, 186bp out of over 12000bp, missed the interacting site between sperm *Bindin* and *EBR1*. *EBR1* is a very complex protein with a variety of domains, each showing different degrees of repetition (Kamei and Glabe 2003). It is possible that these different domains might differ in the degree to which they interact with sperm *Bindin*.

## Results

### *Laboratory assays of pairwise fertilization success under sperm competition*

Fertilization experiments were conducted under a small range in sperm concentrations where typically there is just enough sperm to achieve a high level of fertilization, and most variation in fertilization was caused by differences among crosses and not subtle differences in sperm concentration (Figure 2A). Pairwise reproductive success of males and females under conditions of direct sperm competition in the laboratory revealed a significant effect of male sperm *Bindin* genotype at the two sites in this exon with common variants (variants and frequencies described in Levitan and Stapper 2010); a Serine/Proline substitution (SP) and an insertion-deletion site (InDel)

(Figure 2 B). The same genotypes in females were not significant, but the matching between males and females at these two sites did significantly influence pairwise reproductive success, indicating that the female genotype for sperm *Bindin* predicts success dependent on the degree of matching with the male genotype (Table 1). As noted in field studies, the less common forms of the protein noted in field collections (Proline and the insertion) had higher reproductive success compared to the more common form, with heterozygotes having intermediate success (Levitan and Stapper 2010). This negative frequency-dependent pattern is the predicted outcome when polyspermy is a risk; less compatible egg protein variants can be maintained at higher frequencies because they are less likely to result in polyspermy, these egg protein frequencies then dictate the matching sperm protein variants (Tomaiuolo and Levitan 2010). The novel finding from this laboratory experiment is that males and females matching at either the Serine/Proline or InDel sites had higher reproductive success than partial- or non-matched individuals (Figure 2C). This assortative fertilization is the predicted outcome when sperm are well-mixed and matching sperm and egg recognition loci are associated by linkage disequilibrium (Tomaiuolo and Levitan 2010). Because we demonstrate this matching is based on a specific point substitution and an InDel, we can examine linkage-disequilibrium at these sites with single nucleotide point substitutions found in the egg receptor protein (EBR1).

#### *Diversity, and selection on EBR1 and Bindin*

We sequenced a total of 4277bp of the *EBR1* gene distributed in 8 loci (*tsp1*, *tsp2*, *tsp4*, *cub3*, *cub9*, *hyr4*, *hyr5*, *hyr7*) (Figure 1) in 29 *S. purpuratus* individuals. Out of the 4277bp, 3557bp were sequences without insertions, and 1647bp were coding DNA. In



general, the *Bindin* and *EBRI* genes showed different diversity and recombination patterns, with *EBRI* having different degrees of variation according to the specific repeat and the type of repeat (summarized in Table 2). Haplotype and nucleotide diversity varied within *EBRI* and between *EBRI* and *Bindin*. Haplotypic diversity was high for *Bindin* ( $H=0.83$ ), and variable among the *EBRI* exons. Haplotypic diversity was high for the *hyr* repeats ( $H \geq 0.9$ ), intermediate for the *tsp* repeats and variable and lower for the *cut* repeats (0.28-0.75). These patterns were mirrored in the levels of nucleotide diversity with *hyr*, *tsp* and *cut* exons having the highest, intermediate and lowest diversity respectively. Sperm *Bindin* had intermediate levels of nucleotide diversity compared to these *EBRI* exons.

The three *tsp* repeats analyzed had a combined total of 16 synonymous and 11 nonsynonymous substitutions, out of which 14 of the synonymous and 8 of the nonsynonymous substitutions belonged to *tsp1* and *tsp4*. The two *cut* repeats had three synonymous and eight nonsynonymous substitutions. The three *hyr* repeats had a combined total of 17 synonymous and 44 nonsynonymous substitutions. The three *hyr* repeats' amino acid properties differed at some sites among the three repeats. For example, some nonsynonymous substitutions changed polarity, which affects the hydrophobicity and hydrophilicity of amino acids (Figure 3).

It is very difficult to pinpoint the origin and forces maintaining molecular variation observed at the DNA sequence level from sampling natural populations. There are two basic types of information that we can extract from DNA sequencing data to estimate nucleotide diversity; the number of segregating sites ( $\theta_w$ ), and the mean number of pairwise differences between haplotypes ( $\theta_\pi$ ). These data capture the patterns of

molecular variation on a sample, and are used by tests such as Tajima D' and Fu and Li's D and F to infer selection under the assumption of constant population size (Tajima 1989; Fu and Li 1993). Although rates of offspring recruitment can vary spatially and temporary in *S. purpuratus*, genetic analysis indicate a strong signal that recruits come from a large number of breeding individuals (Flowers et al. 2002). Overall, most exons of *Ebr1* show low levels of nucleotide diversity ( $\theta_\pi$ ) with respect to ( $\theta_w$ ) (Table 2). The negative values of the tests of neutrality over all loci indicate an excess of low frequency polymorphisms, and suggest that in general these proteins are under purifying selection, although we cannot rule out alternative explanations such as population expansion or a recent selective sweep (Table 2).

#### *Linkage disequilibrium between Bindin and EBR1*

The Migrate analysis indicates that *EBR1* shows structure associated with the sperm *Bindin* genotype with respect to the Serine/Proline substitution (Table 3). The *structured* model (individuals binned by their sperm *Bindin* genotype) had a probability of 1.000 whereas the *panmictic* model (no association based on sperm *Bindin* genotype) had a probability of 0.000. This suggests linkage disequilibrium between the two genes (Table 3). In contrast, the Migrate analysis of the presumably neutral microsatellite loci (Table 3) indicated that the *structured* model was not supported (model probability of 0.000), but the *panmictic* model (1.000) was supported. Finding strong evidence for linkage disequilibrium in the interacting GRPs and not with neutral markers provides evidence that the association between sperm *Bindin* and *EBR1* is driven by selection and not by shared demographic history.

We also tested the *panmictic* versus the *structured* model independently at each site to try to determine whether all or some domains of *EBR1* were under linkage disequilibrium with the Serine-Proline site in sperm *Bindin*. The *EBR1* domains associated with sperm *Bindin* Serine /Proline could be involved with the active site where *EBR1* and *Bindin* interact, or could be physically linked genomic locations for this interaction (Table 4). We used a criterion of a ln Bayes Factor < - 2 as evidence of significant support for either the *structured* or *panmictic* model, testing *EBR1* repeats independently with sperm *Bindin* (Table 4). The *structured* model was significantly preferred for *hyr4* (model probability = 1.0), *hyr5* (p = 0.89) and *tsp4* (p = 0.9). In contrast, we found that the *panmictic* model was significantly preferred for *cub9* (p = 0.97), suggesting that particular *EBR1* protein domains vary in the interaction strength with *Bindin*.

To illustrate the relationship between the Migrate results and the classical calculation of LD between two point substitutions, we used the method of Rogers and Huff (Rogers and Huff 2009) to calculate  $R^2$  values based on genotype data for the most common non-synonymous point substitutions in *hyr4* ( $R^2 = 0.07$ ), in which significant LD was noted, and *cub9* ( $R^2 = 0.002$ ), in which no evidence for LD was found (Figure 4). Two regions in *tsp1* had similarly high estimates of LD ( $R^2 = 0.067$ , S.E. 0.013) compared to *hyr4*, while two intronic substitutions surrounding *tsp1* and *hyr4* had low estimates of LD ( $R^2 = 0.006$ , S.E. 0.001) comparable to *cub9*.

Interestingly, repeats in which the population was identified as being *panmictic* were the sites with the lowest nucleotide diversity (Table 2). This suggests that the active site is likely found in *hyr4*, *hyr5*, *hyr7*, *tsp1* or *tsp4*, or that mutations in these regions act

in tandem to influence assortative mating (i.e. by altering folding of the protein).

However, it is important to point out that since we only tested a subset of repeats of all *EBR1* proteins, the actual active site may be located in one of the regions we did not target.

#### *Linkage disequilibrium and recombination within EBR1*

We explored whether there were patterns of linkage within *EBR1* and if considering these patterns provided additional support of either a *structured* or *panmictic* model in relation to sperm *Bindin*. We characterized patterns of linkage within the eight units in the *EBR1* gene and their association with sperm *Bindin*. We compared model fit of the structured model for all linkage combinations of the physically ordered units, *tsp1*, *tsp2*, *tsp4*, *cub3*, *cub9*, *hyr4*, *hyr5*, and *hyr7*. This resulted in 128 patterns from complete linkage of all units (equivalent to concatenation) to completely unlinked units (eight independent loci). The marginal likelihood of the model was used to order the models. We examined these patterns using the Serine-Proline, the InDel site and the most common synonymous substitution (C/T) in sperm *Bindin* as the structuring agent (Table 5). Of the 128 possible linear linkage groups for *EBR1* tested across these three structuring agents, the favored model ( $P = .705$ ) was the *structured* model based on Serine-Proline with all sampled *EBR1* repeats being unlinked. The next most likely model was also structured by Serine-Proline but with *cub3* and *cub9* being linked and all other repeats being unlinked ( $P = 0.295$ ). All remaining models, including all panmictic models and models that subdivided sperm *Bindin* by either the InDel or CT polymorphism, had probabilities approaching zero and LN Bayes factors  $> 40$ . These

results suggest high levels of recombination among the sampled *EBRI* repeats and that the interaction between *EBRI* and sperm *Bindin* likely resides at this Serine-Proline site.

These results indicate a nonrandom association between *Bindin* and different *EBRI* loci (Figure 4). The fact that associations varied between *Bindin* and different *EBRI* repeats (Table 4,5) is not surprising given that minimum recombination rates for all the *EBRI* loci varied between 1 and 7 recombination events (except for *cub3*, for which  $R_m=0$ ) (Table 2). *Bindin* had no recombination events detected.  $R_m$  was consistently higher for the *hyr* repeats than for *cub* and *tsp* repeats (Table 2). The *tsp* repeats had estimated  $R_m$  values ranging from one to three, the *cub* repeats had the lowest  $R_m$ , ranging between zero and three, and *hyr* had the highest  $R_m$ , estimated to be from three to seven (Table 2). Because larger repeat types may be expected to experience more  $R_m$ , we scaled values of recombination ( $R_m'$ ) by dividing by overall repeat size. Scaled results showed the same patterns as the  $R_m$  for each of the target *EBRI* loci (Table 2).

## Discussion

*Strongylocentrotus purpuratus*, like other echinoid species (Palumbi 1999; Levitan and Farrell 2006) exhibits patterns of assortative fertilization under conditions where all males have an equal opportunity to fertilize eggs. In a previous field study, conducted under a range of water flow conditions, population densities, mate distances and gamete release concentrations, features of the sperm *Bindin* protein (the S/P substitution and the InDel) influenced reproductive success (Levitan and Stapper 2010). Here we found that the same two sperm *Bindin* features influenced fertilization in the controlled conditions of the laboratory and further males and females matched at sperm

*Bindin* had higher reproductive success than mismatched individuals. Theory suggests that under conditions of high sperm concentration, when polyspermy is a risk, less compatible egg recognition proteins would be favored and be present in higher frequencies in balancing selection with more compatible proteins. These egg protein variants, with unique compatibility characteristics, select for the maintenance of a matching set of sperm proteins (Tomaiuolo and Levitan 2010). These conditions that result in balancing selection are also predicted to result in linkage disequilibrium between sperm and egg recognition proteins, which is indirectly manifested as the ability to predict female reproductive success by how their sperm *Bindin* genotype matches their male mate's genotype.

We found direct evidence of linkage disequilibrium between sperm *Bindin* and the female protein *EBRI*. We were able to strongly reject the hypothesis that *EBRI* genotypes were independent of sperm *Bindin*, and we were able to reject the idea that this relationship was caused by a shared population history. We are also able to reject that LD is caused by a genomic proximity between these two interacting loci, as they are separated by over a hundred thousand base pairs (Figure 1). Previous work has searched for, and not found, evidence that sperm *Bindin* is expressed in females or their eggs (Nishioka et al. 1990; Minor et al. 1989). As each of these proteins is expressed exclusively on eggs (*EBRI*) or sperm (*Bindin*), the weight of evidence suggests that this pattern is driven by selection during the fertilization process. The fertilization assays suggest that assortative fertilization is the mechanism resulting in LD between these proteins.

Patterns of linkage disequilibrium driven by assortative mating, rather than by genomic proximity, might be a powerful tool for uncovering contemporary patterns of selection. Gamete recognition proteins are known to often, but not always, show the molecular signature of rapid evolution (Palumbi 1999). These signatures include a more rapid divergence across species (Biermann 1998; Hellberg et al. 2000; Swanson and Vaquier 2002; Zigler and Lessios 2003), or a higher degree of polymorphism within species (Metz and Palumbi 1996; Geyer and Pulumbi 2003; Riginos et al. 2006), than predicted by chance. The mechanisms driving these molecular signatures have been debated, and include historic and long-term patterns of reinforcement selection, sexual selection, sexual conflict and pathogen avoidance (Swanson and Vaquier 2002). In contrast, linkage disequilibrium caused by assortative mating has an ephemeral signature. Because these loci are not physically linked by proximity, the degree of linkage disequilibrium is a balance between assortative mating during fertilization and recombination during meiosis. As such, even a single generation of random mating would eliminate this signature; linkage disequilibrium may be an indication of ongoing rather than historic selection. Although this signature of ongoing selection within a population is sensitive to the degree of assortative mating, it might have a lasting evolutionary effect. If selection for assortative mating is strong enough and the conditions for reproductive isolation are met, e.g. (Gavrilets and Waxman 2002), then these associations would become fixed in reproductively isolated populations.

As noted in *S. purpuratus* sperm *Bindin* (Levitan and Stapper 2010) and sea star *EBR1* (Hart 2013), we found an overall pattern of low nucleotide diversity estimated from pairwise differences ( $\theta_\pi$ ) with respect to the number of segregating sites ( $\theta_w$ ), which

could be explained by purifying selection; however there were a few sites of high nucleotide diversity, which could be a result of population contraction, balancing or positive selection. These sites were more common in the *hyr* and *tsp* regions compared to the *cut* regions. In addition, *hyr* and *tsp* regions showed stronger evidence for LD with sperm *Bindin*. Although we only subsampled the full range of repeats in this large and complex protein, it does provide some hints that reproductive compatibility might be associated with the *hyr* and/or the *tsp* regions. Linked sites are likely to be involved in the species and individual compatibility between sperm and eggs, whereas unlinked sites are likely to be involved in protein structure unrelated to compatibility.

Hart (2013) examined *EBR1* in the sea star *Patiria miniata* and noted sites under positive selection but no evidence that the phylogeographic structure noted in neutral markers and sperm *Bindin* were present in the subset of *EBR1* alleles examined. Our study suggests that different regions of *EBR1* can show divergent patterns of association with sperm *Bindin*, and high degrees of recombination within *EBR1*. The failure to find a co-evolutionary signature of sperm *Bindin* with a subset of *EBR1* alleles in *Patiria* might be caused by either their independent evolution or the independence of the subset of *EBR1* regions investigated.

A second interesting finding from Hart's (2013) investigation of the transcriptome of *EBR1* in *Patiria* is that this sea star, like the sea urchin *Strongylocentrotus franciscanus* lack evidence of expression of *hyr* repeats. Our genomic analysis of a single *S. franciscanus* individual indicated the presence of *hyr* repeats. This suggests that regulation of expression of these *hyr* repeats may play a role in patterns of compatibility across species.



We used the program MIGRATE to examine patterns of LD between sperm *Bindin* and *EBRI*. This allowed us to use combinations of unlinked loci to characterize the strength of structure caused by the Serine-Proline polymorphism compared to other polymorphic sites in *Bindin*. In addition, our approach allowed us to use a finite site mutation model instead of the usual infinite site model. The main strength of our approach was the comparison of different degrees of linkage among the repeats within *EBRI*, allowing us to corroborate the qualitative finding of several recombination events between adjacent repeat elements.

In summary, *EBRI* is a large and complex protein that interacts with the sperm *Bindin* protein in sea urchins, influencing reproductive success. Much work needs to be done to characterize this protein, and in particular determine which regions influence patterns of reproductive success both within and across species. In our sampling of this protein we note that regions of the *EBRI* gene are held in LD with sperm *Bindin*, likely via assortative mating. When linkage disequilibrium is caused by assortative mating it provides a powerful tool for examining the relationship between ongoing selection for mating success with the historic signature of selection gleaned from sequence data.

## **Materials and Methods**

### *Laboratory assays of pairwise fertilization success under sperm competition*

We examined if males and females with the same sperm *Bindin* genotype had higher levels of gametic compatibility compared to mismatched mates. On March 12, 2009, patterns of assortative fertilization were examined by independently testing the

eggs of seven females with a competitive assay from the pooled sperm from ten males. Gametes were collected by injecting sea urchins with 0.5 M KCl. Sperm from males were collected “dry” and on ice, while eggs were collected by inverting the female in a dish of seawater. Eggs from each female were diluted to a suspension of ca. 5000 eggs/ml and 1 ml of this suspension was placed into a vial containing 8 ml of filtered seawater (final concentration of 500 eggs/ml). A 0.1 ml aliquot of dry sperm was diluted 1000-fold via serial dilutions, and two ml each of this sperm suspension from all males were mixed in a single vial. A 1 ml aliquot of this mixture was added to each of the vials containing the eggs from each female. A subsample of each sperm suspension was fixed in formalin for later sperm counts.

The sperm dilution used for this experiment was chosen based on prior published experiments (Levitan 1998, 2002 see Fig. 2A) to be in a region where sperm were just high enough to be saturating without resulting in a high degree of polyspermy and where variation in fertilization could be attributed to among cross variation rather than subtle differences in sperm concentration (Fig 2A).

Three hours post fertilization at least 100 eggs from each vial were examined under the microscope for the presence of a raised fertilization membrane or cleavage. The remaining eggs were washed in fresh seawater to remove excess sperm and placed in a jar with 500 ml of filtered seawater. Tube feet samples of all adults were fixed in 95% EtOH for genetic analysis of sperm *Bindin* genotype (as described below). Three days later, 50 larvae from each vial were collected and frozen for paternity analysis using microsatellite loci. We chose to examine reproductive success of adult individuals via genotyping adults and using paternity analysis rather than sequencing larvae for their

gamete recognition proteins, because unpublished data suggest that gametes express the genotype of the adults that producing the gametes rather than gametes expressing their own haploid genotype.

Prior research has indicated that the commonly noted point substitution leading to a Proline for Serine amino acid change and a commonly found InDel result in high to moderate levels of fertilization, while other rarer non-synonymous substitutions result in low levels of fertilization (Levitan and Stapper 2010). Because of our goal to examine genotype matching between parents, adults with these rare genotypes (three males and one female) were excluded from the analysis, leaving seven males and six females in the fertilization experiment. The common alleles (see Levitan and Stapper 2010 for aa code) were Allele A (Serine/Del with a population frequency of 0.58), J (Proline/Del, freq. = 0.16) and K (Serine/Insert, freq. = 0.14). Male genotypes in the crosses were KK, JJ, JK, AJ, AJ, AK, AK; female genotypes were AK, AA, AJ, KK, JK, JK. Both KK/KK and JK/JK crosses (for example) were considered full matches.

For each female, between 35 and 42 larvae were examined for paternity (methods of DNA extraction, PCR and microsatellite loci can be found in Levitan 2008). Pairwise reproductive success was calculated as the product of the proportion of eggs fertilized and the paternity share of each male. Pairwise reproductive success was examined in an ANCOVA testing the main effects of sperm *Bindin* genotype and the covariate of sperm concentration and the degree of genotype matching between the male and female in the pair. To illustrate the fertilization differences between males of different sperm *Bindin* genotypes and males with different degrees of matching with females for their sperm *Bindin* genotype both genotype and matching were treated as main effects (with sperm

concentration as a covariate) to estimate the Least-Square means for each factor (Fig. 2). Least square means adjust the mean value by the other factors in the model including the covariate of sperm concentration.

*DNA extraction and amplification of selected sperm Bindin, EBR1 and microsatellite loci*

Tube feet from 29 individuals from a wild population of *S. purpuratus* were collected from Barkley Sound, British Columbia, for genetic analysis. Individual test tubes containing tube feet from each adult individual in the study were kept at -20°C in ethanol until extraction. Genomic DNA was extracted from tube feet stored in ethanol. Five to seven tube feet per individual were air dried and placed in a 500ml CTAB (2% hexadecyltrimethyl ammonium-bromide) solution with Proteinase K, and incubated overnight at 65°C. DNA was purified from digested tissue by solid phase reversible immobilization (SPRI) using magnetic beads. Extracted DNA was stored at -20°C. A non-neutral locus from *Bindin* exon1 (AF077309) (Biermann 1998), and seven potentially nonneutral loci spanning the *EBR1* gene chosen for the study (*tsp1*, *tsp2*, *tsp4*, *cub3*, *cub9*, *hyr4*, *hyr5*, *hyr7*) were PCR amplified for sequencing (Figure 1). The *Bindin* locus was chosen since it has been previously shown to influence mating success (Levitan and Stapper 2010); the *EBR1* loci were selected from a pilot study conducted to determine potentially informative sites based on levels of polymorphism (data not shown). Locus-specific primers (Table 6) were designed on intronic regions flanking the target exons using the purple urchin genome and the available mRNA sequence (Genbank NM\_214665.1). Thermal cycler programs were optimized for each of the loci. To minimize sequencing errors, a high fidelity polymerase with proofreading capabilities (Invitrogen's Hifi platinum) was used to amplify target regions. Amplification success

was verified via gel electrophoresis. DNA extracted from gonad tissue from one sympatric *S. franciscanus* individual was also amplified at these *EBR1* loci using the above conditions.

The PCR cocktail for each primer pair targeting each specific region consisted of 12.95  $\mu$ l double-distilled water, 2.5  $\mu$ l 10 $\times$  PCR buffer, 1.0  $\mu$ l 2 mM MgCl<sub>2</sub>, 2.5  $\mu$ l 2 mM dNTPs, 0.15  $\mu$ l Platinum HiFiTaq (Invitrogen), 1.2  $\mu$ l 0.5  $\mu$ M forward primer, 1.2  $\mu$ l  $\mu$ M for reverse primer, 1.0  $\mu$ l 10  $\mu$ M bovine serum albumin, and 1.0  $\mu$ l DNA (25 ng/ $\mu$ l). The PCR program was as follows: 5 min at 94°C; 30 cycles of 1 min at 94°C, 30 sec at 56–58°C depending on the primer set (*tsp1*=58°C, *tsp2*=56°C, *tsp3*=58°C, *cub3*=56°C, *cub9*=58°C, *hyr4*=58°C, *hyr5*=58°C, *hyr7*=56°C), and 1 min at 68°C; and 7min at 68°C. Amplification success was verified via gel electrophoresis. In addition, gonad tissue from one sympatric *S. franciscanus* individual was used for DNA extraction, and *EBR1* *hyr* amplification using the above conditions.

#### *Sequencing, sequence alignment, and analysis of Bindin and EBR1*

Successful amplifications of *Bindin* and *EBR1* were selected for cloning. Using a 2.1topo TA vector (Invitrogen), clones were verified via colony PCR using M13 primers, and 6 clones of each locus were selected for sequencing on an Applied Biosystems 3100 Genetic Analyzer using T7 primers.

Sequence alignment was performed with Sequencher, and used to determine the haplotypes of each target region for each of the individuals in the study (sequences deposited with GenBank, Accession numbers KP661601-KP662044). Haplotype sequence data was analyzed with DNAsp and Mega (Librado and Rozas 2009; Tamura et al. 2011). We determined the number of segregating sites, estimated haplotype diversity

(H) and nucleotide diversity ( $\theta_w$  and  $\theta_\pi$ ), and performed tests of neutrality (Tajima's D, and Fu and Li's D and F (Tajima 1989; Fu and Li 1993). We also computed the minimum number of recombination events at each locus using DNAsp (Rozas J and Rozas R 1995). For neutral loci all individuals were genotyped at 7 informative microsatellite loci (loci S4, 4H12, 002 from Cameron et al. 1999; loci Sd52, Sd76, Sd121, Sd156 from Addison and Hart 2002) using Genemapper 4.1 (Applied Biosystems).

#### *Hyr repeat resolution*

Working with repeats of high similarity presents challenges for proper repeat discrimination due to potential amplification of repetitive units. Although most of the *EBR1* gene is repetitive, only the *hyr* repeat region is problematic for repeat resolution since repeat size and sequence similarity allow for occasional amplification of different repeats which could be erroneously identified as haplotypic variants of the same repeat. For this study we only used repeats that could be unambiguously determined via primer specificity, clone size selection, diagnostic sites in intronic sequence regions flanking the exons of interest, and genome mapping. The *hyr4* repeat generated an amplicon 425 bp long, while the *hyr5* repeat was 459 bp long, and the *hyr7* repeat was 713bp long. Both *hyr4* and *hyr5* repeats had the same length of non-coding DNA 5' of the exon (46bp), while the *hyr7* was 359bp long. Besides the longer sequence 5' of the exon, the *hyr7* has a unique sequence 5' of the exon. The 3' end after the *hyr* exon is 136, 170 and 111bp long respectively for *hyr4*, *hyr5* and *hyr7*. The 3' end sequence includes unique sites for the three repeats, and was used to discriminate between *hyr4* and *hyr5* in combination with amplicon length (Table 7).

#### *Analysis of reproductive success*

We analyzed the patterns of fertilization success using a generalized linear model (Table 2). Factors tested were the main effects of serine/proline substitution (SS, SP, PP) and InDel (II, ID, DD), in both males and females, as well as the covariate degree of matching (0, 1, or 2 alleles as a continuous variable) between mates at both the serine/proline and Indel site.

### *Linkage disequilibrium*

Nonrandom mating with respect to GRP genotype can lead to decreased heterozygosity of the GRP loci while leaving neutral loci unaffected. Here we hypothesized that assortative mating due to specific gametic affinities between *Bindin* and *EBR1* genotypes leads to linkage disequilibrium of *Bindin* and *EBR1* in *S. purpuratus*. Different methods for calculating LD are currently available, most commonly  $D'$  and  $R$ ; however, these methods do not take advantage of multilocus marker data (Jorde 2000). In contrast, we are interested in the association of particular alleles of *Bindin* that are known to predict reproduction success (Levitan and Stapper 2010) with all subunits of *EBR1*, combined or individually. We explored two models: (1) the population is *structured* according to a polymorphism in the *Bindin* gene indicating linkage disequilibrium, therefore information on *EBR1* can be inferred from the *Bindin* genotype; and (2) the population is *panmictic* (unstructured), indicating no linkage disequilibrium and assortative fertilization, in which as a result, genotypic information on *Bindin* does not provide information on the genotypic variation on *EBR1*.

We used the program MIGRATE (Beerli 2006; Beerli and Palczewski 2010) to evaluate the two models. We used the Proline/Serine amino acid substitution site in *Bindin* (site 169 in exon 1) which correlates with increased reproduction success of both

males and females (Levitan and Stapper 2010) . We partitioned the sample into Serine/Serine homozygote females and males, and Serine/Proline heterozygote females and males. There were no Proline/Proline homozygote individuals sampled for this examination of linkage disequilibrium.

The two models (*panmictic* versus *structured*) for each sex were tested, individually or in combinations, for the eight presumptively non-neutral *EBRI* loci (*tsp1*, *tsp2*, *tsp4*, *cub3*, *cub9*, *hyr4*, *hyr5*, and *hyr7*) (see Figure 1) and the seven presumptively neutral microsatellite markers (Cameron et al 1999).

In addition to the *Bindin* Serine/Proline site, we also defined the structured model using an InDel of 3 bp to 12 bp length at site 139 in exon 1 in *Bindin* that correlates with reproductive success (Levitan and Stapper 2010).

We established linkage pattern between the subunits of *EBRI* in the light of the potential association with *Bindin* alleles by comparing the *structured* and *panmictic* model for all possible partitionings of *EBRI* in a linear arrangement. This partitioning led to 128 different (physical) linkage groups for which we established a model order using marginal likelihoods.

The MIGRATE run parameters were calibrated so that the settings used for the comparison showed convergence of the Markov chain Monte Carlo sampling method. For the sequence data, we used the following settings: uniform prior distributions for mutation-scaled population size parameters  $\Theta$ , and mutation-scaled migration rates  $M$  over the range of 0.0 to 0.1 and 0.0 to 5000.0, respectively. Four independent chains using different acceptance ratios (temperature settings were 1.0; 1.5; 3.0; 1,000,000.0) were run concurrently. Each chain was a combination of 100 replicates, each of which



discarded the first 10,000 samples as burn-in. A total of 50 million states were visited and 50,000 states were recorded for the generation of posterior distribution histograms. For the microsatellite data we used these settings for each of the 7 loci, except that prior distributions for  $\Theta$  had ranges of 0.0 to 200.0 and for  $M$ , 0.0 to 100.0, and the burn-in was 100,000; a total of 350 million states were visited and 350,000 samples were recorded. The different models were evaluated with marginal likelihoods. These were approximated with the Bézier-quadrature thermodynamic integration as described by Beerli and Palczewski (2010).

We favored the use of MIGRATE to detect linkage disequilibrium between *EBR1* and *Bindin* as a strategy to determine whether the sampled population of *EBR1* and neutral microsatellite loci behaved as one or two populations with respect to *Bindin*, because of the program's ability to successfully use multilocus polymorphic data. Other commonly used methods, such as PAML (Yang 2007), cannot analyze physically unlinked data, for example the different subunits of *EBR1* and *Bindin*. Users of PAML would concatenate all subunits and thus rely on the use of phased genotypic data, in which polymorphic sites in a sequence are grouped to predict possible haplotypes, to make linkage disequilibrium predictions; however, phased data is problematic when considering multiple highly polymorphic loci and the genomic distance between *Bindin* and *EBR1*.

### **Acknowledgements**

We would like to thank Katie Lotterhos and Kim Reuter for assistance with the fertilization assays and Megan Lowenberg, Karalyn Aranow, Maria Wieselmann and

Rebecca Buchwalter for assistance with the molecular work.

## References

- Aagaard JE, Vacquier VD, MacCoss MJ, Swanson WJ. 2010. ZP domain proteins in the abalone egg coat include a paralog of VERL under positive selection that binds lysin and 18-kDa sperm proteins. *Mol Biol Evol.* 27 (1):193-203.
- Addison JA, Hart MW.. 2002. Characterization of microsatellite loci in sea urchins ( spp.). *Mol Ecol Notes* 2.4:493-494.
- Beerli P. 2006. Comparison of Bayesian and maximum likelihood inference of population genetic parameters. *Bioinformatics* 22(3):341-345.
- Beerli, P and Palczewski, M. 2010. Unified framework to evaluate panmixia and migration genetic parameters among multiple sampling locations. *Genetics* 185(1):313-326.
- Biermann CH. 1998. The molecular evolution of *sperm Bindin* in six species of sea urchins (Echinoida: Strongylocentrotidae). *Mol Biol Evol.* 15:1761-1771.
- Callebaut I, Gilgs D, Vigon, Mornon JP. 2000. *HYR*, an extracellular module involved in cellular adhesion and related to the immunoglobulin-like fold. *Protein Sci*9:1382-1390.
- Cameron RA, Minor JE, Nishioka D, Brittan RJ, Davidson EH 1990. Locale and level of bindin mRNA in maturing testis of the sea urchin, *Strongylocentrotus purpuratus*. *Dev. Biol.* 142:44-49.
- Cameron RA, Leahy PS, Britten RJ, Davidson EH. 1999. Microsatellite loci in wild-type and inbred *Strongylocentrotus purpuratus*. *Dev. Biol.* 208:255-264.
- Clark NL, Aagaard JE, Swanson W J. 2006. Evolution of reproductive proteins from animals and plants. *Reproduction* 131.:11-22.

- Clark NL, Gasper J, Sekino M, Springer SA, Aquadro CF, Swanson WJ. 2009. Coevolution of inter-acting fertilization proteins. *PLoS Genet* 5(7): e1000570.doi:10.1371/journal.pgen.1000570
- Dobelli M. 2005. Adaptive speciation when assortative mating is based on female preference for male marker traits. *J Evolution Bio* 18:1587-1600.
- Flowers JM, Schroeter SC, Burton RS 2002. The recruitment sweepstakes has many winners: genetic evidence from the sea urchin *Strongylocentrotus purpuratus*. *Evolution* 56:1445-1453.
- Fu YX, Li WH. 1993. Statistical tests of neutrality of mutations. *Genetics*. 133(3): 693-709.
- Galindo BE, Vaquier VD, Swanson WJ. 2003. Positive selection in the egg receptor for abalone sperm lysine. *P Natl Acad Sci USA*. 100(8): 4639-4643.
- Gao B, Klein LE, Britten, RJ, Davidson EH. 1986. Sequence of mRNA coding for *Bindin*, a species-specific sea urchin sperm protein required for fertilization. *P Natl Acad Sci USA.*, 83(22):8634-8638.
- Gavrilets S, David Waxman. Sympatric speciation by sexual conflict. 2002. *P Natl Acad Sci USA*. 99.16:10533-10538.
- Geyer LB, Palumbi SR. 2003. Reproductive character displacement and the genetics of gamete recognition in tropical sea urchins. *Evolution*. 57.5 1049-1060.
- Hart M. 2013. 2013. Structure and evolution of the sea star egg receptor for sperm *Bindin*. *Mol Ecol*. 22:2143-2156.
- Haygood R. 2004. Sexual conflict and protein polymorphism. *Evolution* 58.7: 1414-1423.

- Hellberg ME, Moy GW, Vacquier VD. 2000. Positive selection and propeptide repeats promote rapid interspecific divergence of a gastropod sperm protein. *Mol Biol Evol.* 17(3):458-466.
- Jorde LB. 2000. Linkage disequilibrium and the search for complex disease genes. *Genome Res.* 10.10: 1435-1444.
- Kamei N, Glabe CG. 2003. The species-specific egg receptor for sea urchin sperm adhesion is *EBRI*, a novel ADAMTS protein. *Genes Dev.* 17(20):2502-7.
- Levitan DR. 2008. Gamete traits influence the variance in reproductive success, the intensity of sexual selection, and the outcome of sexual conflict among congeneric sea urchins. *Evolution.* 62:1305-1316.
- Levitan D R. 2012. Contemporary evolution of sea urchin gamete-recognition proteins: experimental evidence of density-dependent gamete performance predicts shifts in allele frequencies over time. *Evolution.* 66:1722-1736.
- Levitan DR., Ferrell D L. 2006. Selection on gamete recognition proteins depends on sex, density and genotype frequency. *Science.* 312:267-269.
- Levitan DR, Stapper AP. 2010. Simultaneous positive and negative frequency dependent selection on *sperm Bindin*, a gamete recognition protein in the sea urchin *Strongylocentrotus purpuratus*. *Evolution.* 64:785-797.
- Librado P, Rozas J. 2009. DnaSP v5: A software for comprehensive analysis of DNA polymorphism data. *Bioinformatics* 25: 1451-1452 — doi: 10.1093/bioinformatics/btp187.

- Metz EC, Palumbi SR. 1996. Positive selection and sequence rearrangements generate extensive polymorphism in the gamete recognition protein *Bindin*. *Mol Biol Evol.* 13(2), 397-406.
- Metz EC, Robles-Sikisaka R, Vacquier VD. 1998. Nonsynonymous substitution in abalone sperm fertilization genes exceeds substitution in introns and mitochondrial DNA. *Proc Natl Acad Sci USA.* 95:10676-10681.
- Minor J, Gao B, Davidson E. 1989. In: Schatten G, Schatten, H. editors. *The Molecular Biology of Fertilization.*, San Diego(CA): Academic. pp.73-88.
- Moy GW, Springer SA, Adams SL, Swanson W J, Vacquier VD. 2008. Extraordinary intraspecific diversity in oyster *sperm Bindin*. *Proc Natl Acad Sci USA.* 105:1993–1998.
- Nishioka D, Ward RA, Poccia D, Kostacos C, Minor JE. 1990. Localization of *Bindin* expression during sea urchin spermatogenesis. *Mol Reprod Dev.* 27: 181-190.
- Palumbi SR. 1999. All males are not created equal: fertility differences depend on gamete recognition polymorphisms in sea urchins. *Proc Natl Acad Sci USA.* 96:12632-12637
- Payne RJH, and Krakauer DC. 1997. Sexual selection, space and speciation. *Evolution* 51:1-9
- Pennington T. 1985. The ecology of fertilization of echinoid eggs: the consequences of sperm dilution, adult aggregation, and synchronous spawning. *Biol Bull.* 169:417-430
- Pujolar JM, Pogson G H. 2011. Positive Darwinian selection in gamete recognition proteins of *Strongylocentrotus* sea urchins. *Mol Ecol.* 20(23)2011:4968-4982.

- Riginos C, Wang D, Abrams AJ. 2006. Geographic variation and positive selection on M7 lysin, an acrosomal sperm protein in mussels (*Mytilus* spp.). *Mol Biol Evol.* 23.10 : 1952-1965.
- Rogers AR, Huff C. 2009. Linkage disequilibrium between loci with unknown phase. *Genetics*, 182(3):839-844.
- Rozas J, Rozas R. 1995. DnaSP, DNA sequence polymorphism: an interactive program for estimating Population Genetics parameters from DNA sequence data. *Comput Appl Biosci.* 11:621-625.
- Swanson WJ, Vacquier VD. 1998. Concerted evolution in an egg receptor for a rapidly evolving abalone sperm protein. *Science* 281.5377: 710-712.
- Swanson WJ, Vaquier VD. 2002. The rapid evolution of reproductive proteins. *Nat Rev Genet* 3(2):137-144.
- Tajima F. 1989. Statistical method for testing the neutral mutation hypothesis by DNA polymorphism. *Genetics.* 123(3):585-595.
- Tamura K, Peterson D, Peterson N, Stecher G, Nei M, Kumar S. 2011. MEGA5: Molecular Evolutionary Genetics Analysis using maximum likelihood, evolutionary distance, and maximum parsimony methods. *Mol Biol Evol* 28:2731-2739.
- Tomaiuolo M, Levitan DR. 2010. Modeling how reproductive ecology can drive protein diversification and result in linkage-disequilibrium between sperm and egg proteins. *Amer Nat* 176:12-25.
- Turner LM, Hoekstra HE. 2008. Causes and consequences of the evolution of reproductive proteins. *Int Devel Biol.* 52(5): 769.

- Vaquier VD, Moy GW. 1977. Isolation of *Bindin*: the protein responsible for adhesion of sperm to sea urchin eggs. *Proc Nat Acad Sci USA*. 74(6):2456-2460.
- Yang, Z. 2007. PAML.4: a program package for phylogenomic analysis by maximum likelihood. *Mol. Biol. Evol.* 24:1586-1591.
- Zigler KS, Lessios HA. 2003. 250 million years of *Bindin* evolution. *Biol Bull* 205(1): 8-15.
- Zigler, KS, McCartney, MA, Lessios HA, Levitan DR 2005. Sea urchin *bindin* divergence predicts gamete compatibility. *Evolution* 59:2399-2404.



## Figure Legends

Figure 1. *Bindin* and *EBRI* loci used in this study showing genomic distances between adjacent sequenced regions.. The genomic scaffold annotations of *Bindin* (Scaffold 66693:2682-4896) and *EBRI* target loci (scaffold 81713:372025-442027) suggest that the two genes are not physically linked. The minimum distances between *EBRI* and *Bindin* are 127760 bp if scaffold 66693 (*Bindin*) is the continuing scaffold immediately following the end of 81713 (*EBRI*), and 382490bp if scaffold 81713 follows scaffold 66693.

Figure 2. Fertilization success. (A) Proportion of eggs fertilized from crosses in study and from 29 random crosses from the population in Barkley Sound as a function of sperm concentration. Experiments were conducted in the region where slight variation in sperm concentration explained little of the variation in fertilization (B) Reproductive success of individuals (Least Squared Means of Pairwise Reproductive Success) with different forms of the *Bindin* protein Proline (P)/Serine (S) substitution and the presence (I) or absence (D) of an insertion in either homozygous or heterozygous males. (C) Reproductive success of males and females according to their level of matching at either the Serine/Proline or the InDel site. Numbers indicate samples sizes for each genotype. Least squared means adjusts the mean value by other factors in the statistical model, including sperm concentration (Table 1).

Figure 3. Heatmap of hydrophobic (light gray) and hydrophilic (dark gray) sites in aligned *Hyr* repeats. The three most common haplotypes for each repeat type are shown.

Some amino acid differences altering the hydrophobicity of the domain are repeat-specific.

Figure 4. Frequencies of Serine homozygous and Serine-Proline heterozygous genotypes as a function of *EBR1* common point substitutions. In *cub9* (bp 369) *Bindin* and *EBR1* genotypes are independent (*panmitic* model favored, Table 5), while in *hyr4* (bp 101) significant LD was noted (*structured* model favored, Table 5) and there is a deficiency of SS genotypes associated with AA individuals and an excess of SP genotypes with GG individuals.  $R^2$  values calculated using genotype data using Rogers and Huff (2009).

Figure 1

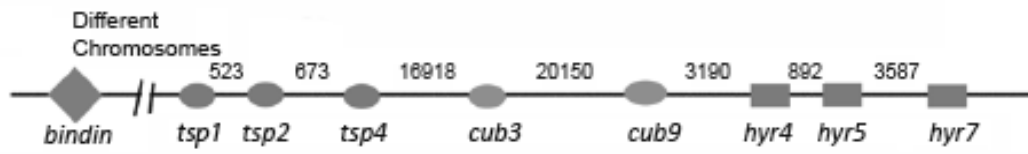


Figure 2

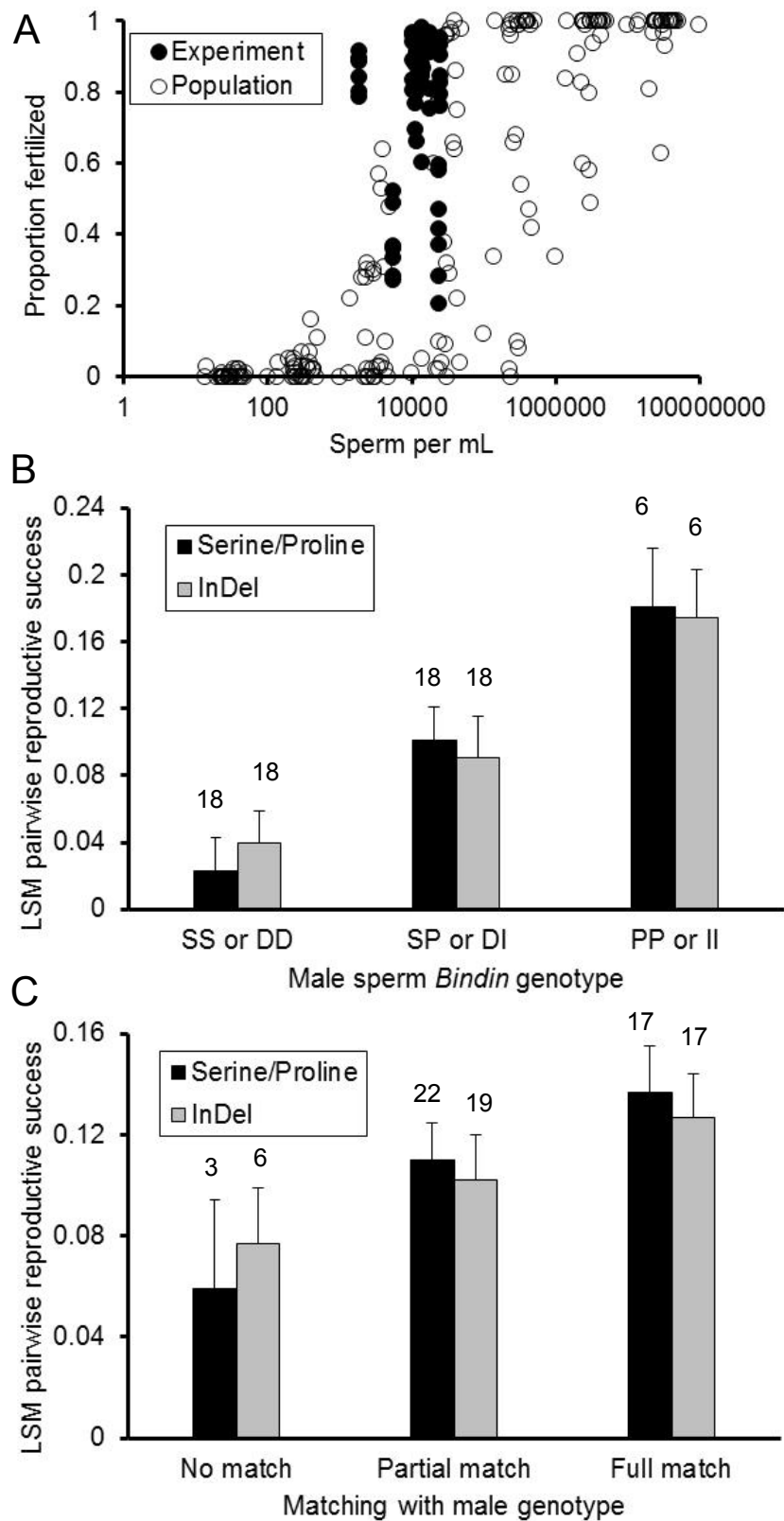


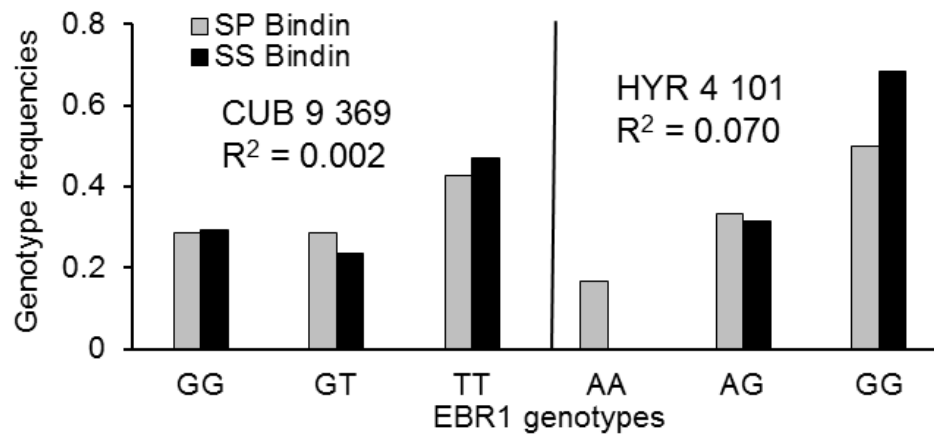
Figure 3

5'	D	N	E	I	P	V	F	S	G	C	P	S	D	Q	N	V	T	T	D	I	G	N	A	T	A	V	V	I	W	T	P	P	
Hyr4	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.
Hyr4	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	N	.	.	.	.	.	.	.	.	.	.	.	.	.	
Hyr4	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	
Hyr5	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	
Hyr5	.	.	K	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	T	.	.	.	
Hyr5	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	
Hyr7	.	.	.	.	.	.	I	.	.	.	.	.	.	.	.	.	.	A	.	.	.	.	.	.	.	.	.	.	T	.	.	.	
Hyr7	.	.	.	.	.	.	I	.	.	.	.	.	.	.	.	.	.	A	.	.	.	.	.	.	.	.	.	.	T	.	.	.	
Hyr7	.	.	.	.	.	.	I	.	.	.	.	.	.	.	.	.	A	.	S	.	.	.	.	.	.	.	.	.	T	.	.	.	

	T	A	T	D	N	S	G	N	L	T	L	T	S	T	N	N	P	G	D	D	F	P	I	G	N	N	T	V	T	Y	S	A	
Hyr4	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.
Hyr4	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.
Hyr4	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	N	.	.	.	.	.	.	.	.	.	.	.	.	
Hyr5	.	.	.	.	.	.	.	S	Q	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	
Hyr5	.	.	.	.	.	.	.	S	Q	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	
Hyr5	.	.	.	.	.	.	.	S	Q	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	
Hyr7	.	.	.	.	.	.	V	Q	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	
Hyr7	.	.	.	.	.	.	V	Q	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	
Hyr7	.	.	.	.	.	.	V	Q	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	

	S	D	Y	A	G	N	T	E	T	C	T	F	F	V	V	V	S	3'	Freq
Hyr4	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	0.423
Hyr4	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	0.135
Hyr4	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	0.077
Hyr5	.	.	D	.	.	.	.	.	.	.	.	.	.	.	.	I	.	.	0.615
Hyr5	.	.	D	.	.	.	.	.	.	.	.	.	.	.	.	I	.	.	0.096
Hyr5	.	.	D	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	0.058
Hyr7	N	.	D	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	0.344
Hyr7	N	.	D	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	0.156
Hyr7	N	.	D	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	0.094

Figure 4



## Tables

Table 1. General Linear Model of influence of sperm *Bindin* genotype on pairwise reproduction success in the laboratory fertilization assay. Factors tested were the main effects of Serine/Proline substitution (SS, SP, PP) and Indel (II, ID, DD) in both males and females, as well as the covariate of sperm concentration and the degree of matching (0, 1, or 2 alleles as a continuous variable) between mates at both the serine/proline and Indel site.

Source	DF	Type III SS	MS	F	Prob. > F
Male S/P	2	0.0219	0.0109	6.27	0.0052
Male Indel	2	0.0308	0.0154	8.81	0.0009
Female S/P	1	0.0011	0.0011	0.64	0.431
Female Indel	2	0.0034	0.0017	0.97	0.3903
Match S/P	1	0.009	0.009	5.16	0.0302
Match Indel	1	0.0086	0.0086	4.93	0.0338
Sperm	1	0.0015	0.0015	0.86	0.3612
Error	31	0.05408	0.0017		
Total	41	0.1079			

Table 2. N=number of sequences, nt=nucleotide (excluding incomplete data and gaps, indels), cds nt=coding sequence, S= segregating sites, Syn=synonymous substitutions, NSyn=nonsynonymous substitutions, H= Haplotypic diversity,  $\theta_w$  = nucleotide diversity estimated from number of segregating sites,  $\theta_\pi$  = nucleotide diversity estimated from mean number of pairwise differences, D(t)= tajima D, D(f) Fu and Li's D, F(f)= Fu and Li's F, Rm= minimum number of recombination events. Rm' = minimum number of recombination events scaled by locus size (Rm/nt). Statistical significance Tajima D=\*0.05,\*\*0.01. Statistical significance Fu and Li D and F=\*0.05,\*\*0.02.

Locus	N	nt	cds nt	S	Syn	NSyn	H	$\theta_w$	$\theta_\pi$	D(t)	D(f)	F(f)	Rm	Rm'
<i>tsp1</i>	54	369	198	33	7	6	0.84	0.0196	0.0087	-1.851*	-3.212*	-3.240**	3	0.008
<i>tsp2</i>	54	513	192	15	2	3	0.718	0.0064	0.0025	-1.939*	-2.876*	-3.027*	1	0.001
<i>tsp4</i>	46	456	177	26	7	2	0.890	0.0130	0.0052	-2.099*	-3.163*	-3.315**	3	0.006
<i>cub3</i>	46	248	168	5	1	4	0.283	0.0046	0.0012	-2.080*	-3.050*	-3.194*	0	0
<i>cub9</i>	50	385	180	10	2	4	0.756	0.0058	0.0047	-0.542	0.1710	-0.075	3	0.008
<i>hyr4</i>	52	419	242	56	9	20	0.965	0.0301	0.0141	-1.795	-2.544*	-2.707*	5	0.012
<i>hyr5</i>	52	459	242	41	7	12	0.905	0.0212	0.0072	-2.243**	-2.772*	-3.076*	7	0.015
<i>hyr7</i>	32	708	242	55	1	12	0.981	0.0193	0.0078	-2.268**	-2.988*	-3.250*	7	0.009
<i>bindin</i>	58	237	252	13	6	7	0.836	0.0188	0.0063	-1.358	-1.135	-1.434	0	0



Table 3. Comparison of a panmictic (model 1) and a structured model (Model 2) for 8 loci in *EBR1* and 7 microsatellite markers. For the structured model, individuals were partitioned using the serine/proline polymorphism in sperm *Bindin*. Male and female sea urchins were analyzed independently. Ln Bayes factor (LnBF) was calculated as the difference of the logarithms of the marginal likelihood (ln ml) of model 1 and model 2.

Data	Sex	Ln ml		LnBF	Probability of model 2
		Model 1	Model 2		
ebr1	Female	-4531.25	-4517.49	13.76	1.0000
	Male	-5036.72	-5022.52	14.20	1.0000
microsatellite	Female	-1255.38	-2182.43	-927.05	0.0000
	Male	-2739.81	-4854.05	-2114.24	0.0000

Table 4. Comparison of a panmictic (Model 1) and a structured (Model 2) for 8 loci in *EBRI*. For the structured model, individuals were partitioned using the serine/proline polymorphism in *Bindin*. Ln Bayes factor (Ln BF) was calculated as the difference of the logarithms of the marginal likelihood (Ln ml) of model 1 and model 2. A model was deemed to be significantly favored when the difference in Ln BF < -2. The structured model was favored in three loci (*Tsp 4*, *Hyr 4*, *Hyr 5*), the panmictic model was favored in one locus (*Cub 9*).

<b>Ebr1 Locus</b>	<b>Model</b>	<b>Ln ml</b>	<b>Ln BF</b>	<b>Model probability</b>
Tsp 1	<b>2</b>	-874.77	0	0.7974
	1	-876.14	-1.37	0.2026
Tsp 2	<b>1</b>	-917.57	0	0.8699
	2	-919.47	-1.9	0.1301
Tsp 4	<b>2</b>	-933.34	0	0.8984
	1	-935.52	-2.18	0.1016
Cub 3	<b>2</b>	-705.93	0	0.6985
	1	-706.77	-0.84	0.3015
Cub 9	<b>1</b>	-1406.43	0	0.968
	2	-1409.84	-3.41	0.032
Hyr 4	<b>2</b>	-1246.24	0	1
	1	-1271.51	-25.27	0
Hyr 5	<b>2</b>	-1092.26	0	0.8938
	1	-1094.39	-2.13	0.1062
Hyr 7	<b>2</b>	-1535.72	0	0.8115
	1	-1537.18	-1.46	0.1885

Table 5. Comparison of a panmictic (Model 1) and a structured model (Model 2) for 8 loci in *EBR1* using the fertilization informative mutations of the Serine/Proline substitution (SP) and an insertion/deletion (InDel). For the structured model, individuals were partitioned using the serine/proline polymorphism, and InDel polymorphisms in *Bindin*. Different levels of *EBR1* linkage were also tested for an association with *Bindin* and here we report the two best models for *EBR1* linkage. The structured model was favored for the Serine/Proline substitution (panmictic models not reported as they all had model probabilities  $< 0.0000$ ), the panmictic model was favored for the InDel. Ln Bayes factor was calculated as the difference of the logarithms of the marginal likelihood (Ln ml) of model 1 and model 2. For details see Methods.

Data	Model	Ln ml	Ln BF	Model probability
SP	2(All Ebr1 loci unlinked)	-8731.2	0	0.7047
	2(Cub3 and Cub9 linked)	-8732.1	-0.87	0.2953
InDel	1(Cub3 and Cub9 linked)	-8774.9	-43.68	0
	1(All Ebr1 loci unlinked)	-8782.9	-51.22	0
	2(Cub3 and Cub9 linked)	-8803.9	-72.62	0

Table 6. Genomic regions targeted for sequencing and primers.

Locus	Sequenced genomic location	Primers	
Bindin	scaffold66693: 2975-3210	f'cttcatctcggggcattctc	r'ttgggtgactacagcgtga
Tsp1	scaffold81713: 382658-383026	f'agaccagtatcggaacacc	r'atcagattcattgacgcaca
Tsp2	scaffold81713: 383549-384061	f'cgaaatcagcagcctagc	r'tgtgaggaggcaggtctttg
Tsp4	scaffold81713: 384734-385272	f'gctatactgggtgcttcttc	r'ccttgattgaagccacgag
Cub3	scaffold81713: 402190-402630	f'cagggtgaaatgtgaggaac	r'tccgatttcagagcatttc
Cub9	scaffold81713: 422780-423414	f'gacctgctcttggggtatga	r'gctacggtcacaaagggttcg
Hyr4	scaffold81713: 426604-427027	f'atggggctcgttatgtatgg	r'cactgggcatccactaaaga
Hyr5	scaffold81713: 427919-428377	f'acagggtcattatgtttgc	r'caaagaaggtgcatatcggttc
Hyr7	scaffold81713: 431964-432676	f'ccatcaccactaccgatt	r'tgctgtaggtggtgtccaag

Table 7. *Hyr* repeat determination. Each target exon was confirmed by having an amplicon of a specific size with specific flanking non coding sequences (ncds).

	Length of per	Length of 5' ncd	Length of 3' ncd
<i>hyr4</i>	425	46	136
<i>hyr5</i>	459	46	170
<i>hyr7</i>	713	359	111

Sequence alignment 5' upstream of exon

*hyr4* **GCGCAGTCA**agttccatctagtt**AtatCtcat**--aatcttttcaca

*hyr5* **GCGCAGTCA**agttccatctagtt**AtatCtcat**--aatcttttcaca

*hyr7* **TCATTTCGA**agttccatctagtt**CtatTtcatATA**aatcttttcaca

Sequence alignment 3' downstream of exon

*hyr4* ggtaag**TAG**ttctcttcta**CtcaT**cttctgaaaaacagtttgta**CtctagtgagagctggaatcagattttagaaatgtctacattttC**

*hyr5* ggtaag**TAC**ttctcttcta**CtcaA**cttctgaaaaacagtttgta**GtctagtgagagctggaatcagattttagaaatgtctacattttA**

*hyr7* ggtaag**AGC**ttctcttcta**AttcaT**cttctgaaaaacagtttgta**CtctagtgagagctggaatcagattttagaaatgtctacattttC**

**Financial Disclosure**

This work was funded by the National Science Foundation DEB 0822626. The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.