

Response to Lior Pachter's criticism of our RECOMB 2007 paper.

In [RECOMB 2007](#), we published a paper titled: "[Network motif discovery using subgraph enumeration and symmetry-breaking](#)" by Grochow and Kellis. In his blog post dated 2/12/14, Lior Pachter reveals that he was one of the reviewers of our paper, and claims that the paper is incorrect and trivial.

We respectfully disagree, and we refute each of his criticisms below:

Criticism 1 (subjective): Pachter claims our results are trivial.

Our algorithm advanced the state of the art in motif finding, enabling the identification of larger network motifs than ever before. It achieved that by introducing several key ideas, which have now been built upon in multiple algorithms to identify even larger motifs:

- 1) We reversed the traditional directionality of the network motif discovery problem, to start with subgraphs and then search for them in the network, rather than scanning the network and enumerating all subgraphs of a given size.
- 2) We provided a solution to the subgraph isomorphism and enumeration problem that runs exponentially faster on real-world networks, by recognizing self-symmetries of network motifs and imposing symmetry-breaking conditions that eliminated repeated instances.
- 3) In addition, using our algorithm we discovered the role of combinatorial effects in larger motifs, and introduced a motif clustering score and computed the clustering properties of all discovered motifs. We explicitly described these both in our original submission and in our published paper.

We stand by the novelty and contribution of our work, which has not only stood the test of time since 2007 (see our response to Criticism 6 below), but has since been replicated and built upon, and is now part of the cited literature in the field of network motif research.

Prior to our algorithm, no one could find *even a single* motif of 15 nodes, let alone dream of looking for motifs of 20 nodes, 30 nodes, or more (we successfully searched for 31-node subgraphs, for example, although we didn't find any of statistical significance). Even today, only two algorithms are able to reach these motif sizes: (1) [Omidi, Schreiber, Masoudi-Nejad. MODA. Genes Genet. Syst. 2009 Oct, 84(5):385-95]; and (2) [Ribeiro and Silva. G-Tries: an efficient data structure for discovering network motifs. ACM Symposium on Applied Computing (SAC '10)]. Both of these algorithms build on our work, and would not be able to reach such motif sizes without the ideas first introduced in our paper.

Criticism 2 (technical): Pachter claims that the 15-node motif in Fig. 3 is solely due to combinatorial effects due to the abundance of cliques relative to Erdős-Renyi random graphs.

As explicitly stated in our paper, the enrichment of the motif in Fig. 3 is due to combinatorial effects. However, it is not as simple as Pachter suggests in his blog:

- 1) First, **Pachter's claim** that the high frequency of the 15-node graph of Fig. 3 is merely due to the presence of a large clique and independent set **is mathematically incorrect**. It is due to the presence of a vertex v that is connected to *every node* in a large independent set, such that none of the nodes in the independent set are connected to *any* of the other 14 nodes in the graph on the right of the figure (and similar descriptions for other nodes in the graph). Even in a random graph model that was somehow rigged to have large cliques and large independent sets, it would be highly unlikely to have many independent sets with that property.

- 2) Second, Pachter states that we define motifs relative to subgraph frequency in an Erdős-Renyi model, which is factually incorrect, and a key to our next point (3) below. Erdős-Renyi random graphs exhibit with very high probability a Poisson degree distribution, which is very different from the long-tailed degree distribution observed in most real-world networks, including those examined in our paper. For that reason, we used a standard Markov chain Monte Carlo edge-swapping technique in our null model, similar to nearly all network motif papers. In order to preserve not just the degree distribution but also the distribution of 3-node motifs, we combined the MCMC edge-swapping with a simulated annealing technique.
- 3) Pachter's criticism hinges on the assumption that both 9-node cliques and 12-node independent sets are over-represented in the real network relative to our control network. However, **this is incorrect**. Although cliques would be over-represented in our network relative to an Erdős-Renyi control, the data show that cliques are neither over- nor under-represented in the real network relative to the control we actually use.

Thus, **Pachter's specific criticism about Fig. 3**, on which hinges his criticism of combinatorial effects, **is incorrect**.

More generally, **not all large motifs arise from combinatorial clustering effects**, and our clustering score helps distinguish large motifs that arise in a combinatorial manner from those that do not. Biologically, the discovery of combinatorial effects for some of the large network motifs has taught us something new about the nature of some large network motifs and the organization of large networks. Mathematically, the presence of combinatorial effects doesn't mean our algorithm is flawed. On the contrary, it shows that the definition of network motifs needs to incorporate clustering properties, which was directly addressed in our paper, both in the original submission and in the published form.

These properties, which in retrospect make a lot of sense, had not been previously reported, and were to our knowledge first recognized as a direct result of our work, as no other algorithms had managed to recognize large network motifs and identify all their instances.

Criticism 3 (technical): Pachter claims our algorithm is essentially a heuristic for finding cliques and stable sets.

Pachter claims that "the Grochow-Kellis method is essentially a heuristic for finding cliques and stable sets in graphs." This is a **gross misrepresentation** of our algorithm, both at a theoretical level and at the level of specific examples:

- 1) Pachter claims to have inferred this claim based on his interpretation of the example of Fig. 3, which we refuted in the previous section. Moreover, even if combinatorial effects accounted for the motif enrichment in that figure, **it is logically incorrect** to infer from this one example that our algorithm merely finds cliques and stable sets.
- 2) As Pachter himself points out, the combinatorial effects present in Fig. 3 would also be present not just near cliques or independent sets, but near any subgraph H that itself contained many isomorphic subgraphs, **which our algorithm would also find**. There are many possibilities for such an H: notably any vertex-transitive graph, or any graph that differs from a vertex-transitive graph in only a few edges. The class of vertex-transitive graphs is quite large (it includes all Cayley graphs of groups), and the presence or absence of such graphs in real-world networks would beg for a theoretical explanation.
- 3) A different example, in Figure 3-6 from Joshua Grochow's thesis (from August 10, 2006, a month before we even submitted to RECOMB), is an *anti*-motif of 10-nodes with 95,754 instances covering 472 nodes, that is clearly not due to the presence of cliques or independent sets. Although the instances of this 10-node graph are heavily clustered, they are much less clustered than the example in Fig. 3. Furthermore, as an *anti*-motif, it appears *less* frequently in the real network than in the null model, and hence its clustering and significance in

the real network cannot be due merely to the presence of cliques or independent sets in the real network compared to the null model.

Criticism 4 (review process): Pachter claims that we got the idea for clustering and the explanation of 27,720 from his review.

Pachter spends three paragraphs explaining how he “took many courses on enumerative and algebraic combinatorics” and “there are some numbers that combinatorialists just know”. He continues: “When I saw the number 27,720 I had a hunch”, but gives us credit that “**the idea may have entered my mind because of the picture of the motif**”. And then concludes, four paragraphs later that “Grochow and Kellis [...] added a section titled ‘The role of combinatorial effects’, in which they explained the origins of the number 27,720 that they gleaned from my report”.

Even if we had gotten these ideas from his report, it is expected that the authors follow the reviewer’s recommendations when revising their paper; although sometimes this is acknowledged by thanking an anonymous reviewer, this is not universally practiced and in some fields even discouraged. But in this case, we didn’t get these ideas from the referee report:

- 1) First, we had explicitly drawn the motif in a way that reveals the combinatorial nature of the cliques, which Pachter semi-acknowledges, saying “the idea *may* have entered my mind because of the picture of the motif”.
- 2) Second, what Pachter claims was ‘gleaned’ from his report in the revised version of our paper was in fact included in Joshua Grochow’s thesis which was submitted on August 10, 2006 to MIT, 1 month before we even submitted the paper to RECOMB and 4 months before we received Pachter’s review. Months before Pachter’s review of our paper, Josh’s thesis wrote: “the frequency of the motif is due to choosing 3 nodes from the large independent set on the left and four from the large clique at the bottom of the figure,” and includes the exact same equation “ $(12 \text{ choose } 3)(9 \text{ choose } 4)=27,720$ ” (Figure 3-5 legend on page 50), while the number or the equation didn't even appear in Pachter's critique.

In any case, it should be clear to the reader that we knew of these ideas prior to receiving any reviews. Moreover, the suggestions that somehow we had not recognized the combinatorial nature of these examples, that it was news to us, that we had to somehow deal with them, and that his review was somehow a revelation to us are **clearly refuted by the discussion in Section 3.4.1 of [Joshua Grochow’s thesis](#)**, finalized and published months before we received any reviews and one month before we even submitted our paper to RECOMB.

As for the more general recognition that large network motifs cluster due to combinatorial effects, that very property would have remained unrecognized without our paper. The **role of clustering only became apparent thanks to our algorithm**, which could for the first time find all instances of such large motifs. Moreover, we had already developed and applied a motif clustering score to specifically address that point ever since our original submission. In responding to the reviews, we expanded our discussion of these ideas and included the section titled “The role of combinatorial effects”, which was directly responsive to the critiques, as is expected for RECOMB and any other scientific conference or journal.

Criticism 5 (technical): Pachter claims our algorithm is incorrect and doesn’t achieve its stated goal.

The stated goal of our algorithm is to find all instances of a query graph H in a graph G, and **our algorithm correctly achieves this stated goal**, as a simple proof of correctness shows. This can then be used to evaluate the statistical significance of the frequency of occurrences of H in G relative to various null models, according to the

standard definition of network motifs. We applied this to discover network motifs, but many other applications are enabled by our algorithm, such as the study of motif clustering, which we demonstrated in our paper.

For small graphs H up to 8 nodes, our algorithm can also be used to *exhaustively* evaluate all subgraphs of that size and their significance (this is larger than anything that had been achieved ever before). For larger sizes, we do not claim that we can find *all* motifs (the sheer number of possible subgraphs simply prohibits this), nor do we claim to find the most significant ones (given the enormous search space, even our sampling approach is not guaranteed to encounter all over-represented subgraphs).

Pachter claims that our algorithm is incorrect because it found a 15-node motif but not a more significant 18-node motif which can be formed combinatorially from the same nodes that make up our 15-node motif. This criticism is analogous to blaming the tallest building in the world for failing to add a few floors. Before our work, no known algorithm could even discover motifs of size 9, let alone 15, 20, or more. We **surpassed the state of the art**, introduced new methods and key ideas (reversing the search, symmetry breaking) that enabled multiple new algorithms since our work, but we did not claim any guarantees about finding the largest motifs, nor the most significant motifs. We can exhaustively search all 8-node motifs, but we use sampling to find larger motifs and anti-motifs, and thus it should be clear that we do not guarantee discovery of all larger motifs.

Criticism 6 (subjective): Pachter claims our solutions are uninteresting and that we should retract our work.

Pachter compares our algorithm to the infamous arsenic paper and states “When you discover your work is flawed, the correct response is to retract it”. However, despite his technical criticisms, our work is neither flawed nor trivial, as we discussed above. Moreover, in contrast to the community uproar caused by the arsenic-based-life-form paper, our algorithm has stood the test of time, has been **replicated and extended, and forms the foundation of much of the subsequent work on discovery of large network motifs.**

Since the conference was 7 years ago, we can look back at how influential each of the papers has been, using a simple metric of citations (we recognize this is an imperfect metric, but it is usually an okay gauge at a high level). From all 38 papers published in RECOMB 2007, the typical paper has 14 citations, while ours has 75 citations, and is **the 2nd most cited paper of the conference** (after the work of Bonnie Berger’s lab who received the RECOMB test-of-time award in 2010). This indicates *significant interest in our work* that continues to this date (the latest citation was ironically on February 14, two days after Pachter’s blog post), and a great impact on the community, which would have been lost if the positive reviews of our paper had not outweighed the negative score given by Pachter.

In conclusion, we have demonstrated that, in contrast to Pachter’s criticisms, **our paper is correct, non-trivial, impactful, and has stood the test of time.**

Joshua Grochow and Manolis Kellis
joshuag@alum.mit.edu, manoli@mit.edu
February 16, 2014