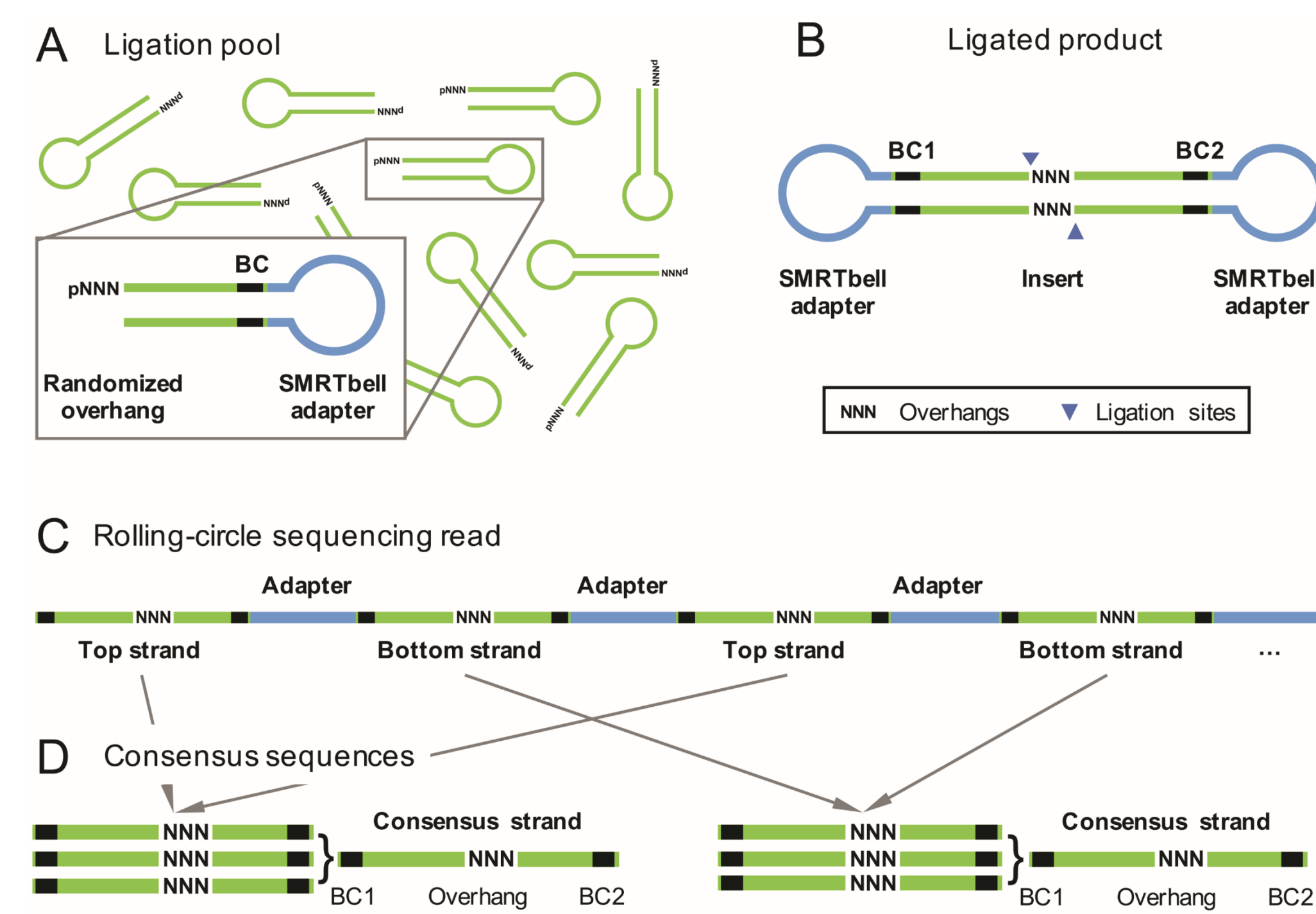


Vladimir Potapov¹, Jennifer L. Ong¹, Rebecca B. Kucera¹, Bradley W. Langhorst¹, Katharina Bilotti¹, John M. Pryor¹, Eric J. Cantor¹, Barry Canton², Thomas F. Knight², Thomas C. Evans, Jr.¹, and Gregory J. S. Lohman¹
¹New England Biolabs, Ipswich, MA 01938, USA ²Ginkgo Bioworks, Boston, MA, 02210, USA

Abstract

The one-pot assembly of long DNA sequences from multiple component parts is key to the rapid generation of constructs for modern synthetic biology. Methods for the one-pot assembly of multiple fragments linked by short overhangs (e.g. Golden Gate) depend on accurate and unbiased ligation. Design of junctions to date largely depends on the use of rules of thumb and empirical success, rather than detailed data on ligase fidelity and bias. In this study, we have applied Pacific Biosciences Single-Molecule Real-Time sequencing technology to directly measure the ligation frequency of every possible 5'-four-base overhang pairing in a single experiment. This comprehensive data set has been applied to predict the accuracy of Golden Gate assembly (GGA) using the Type IIS restriction enzyme BsaI. Ten fragment assemblies were designed based on the ligation data with junctions predicted to result in high or low fidelity assembly. Experimental results confirmed not only the overall accuracy, but the specific mismatch ligation errors observed and their relative frequency. The data was further used to design 12- or 24-fragment assemblies of the lac operon, which were shown to assemble with high fidelity and efficiency. Thus, ligase fidelity data allows the prediction of high-accuracy overhang pair sets with greater flexibility in design than the rules of thumb, allowing assembly of >20 fragments at high-accuracy junction points even within defined coding regions without modification of the native DNA sequence.

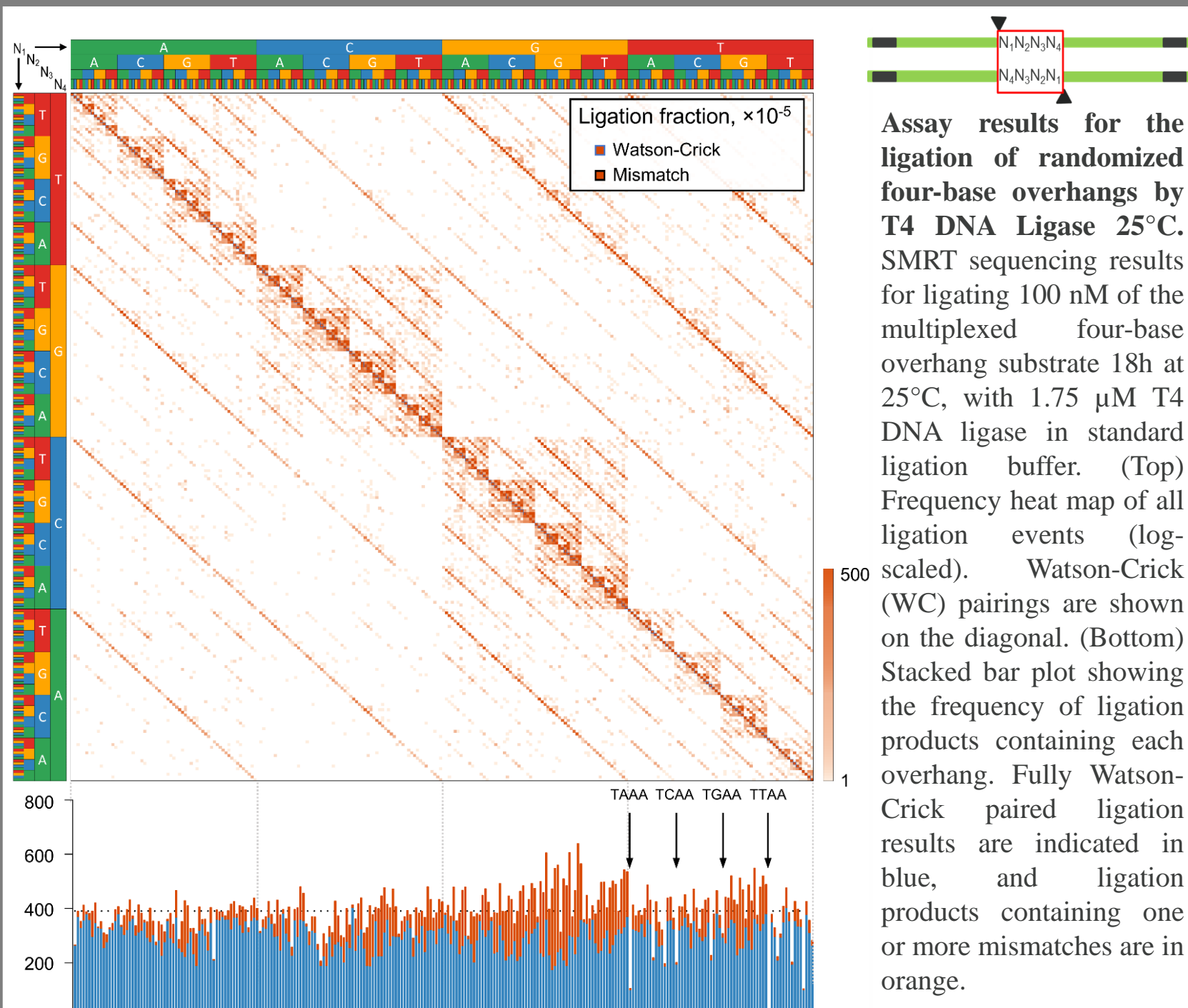
Schematic of Multiplexed Ligation Profiling Assay



(A) Libraries containing randomized three-base overhangs were synthesized and ligated with T4 DNA ligase under various conditions. The hairpin substrates contain the SMRTbell adapter sequence as well as an internal 6-base random barcode used to confirm strand identity and monitor the substrate sequence bias derived from oligonucleotide synthesis. (B) Ligated substrates form circular molecules, in which a double-stranded insert DNA is capped with SMRTbell adapters. (C) Ligated products were sequenced utilizing PacBio SMRT sequencing, which produced long rolling-circle sequencing reads. Sequencing reads are comprised of regions corresponding to top and bottom strands separated by regions corresponding to SMRTbell adapters. (D) Consensus sequences were built for the top and bottom strands independently, allowing extraction of the overhang identity and barcode sequence.

Methods/Results

Fidelity and Bias: Four-Base Overhangs



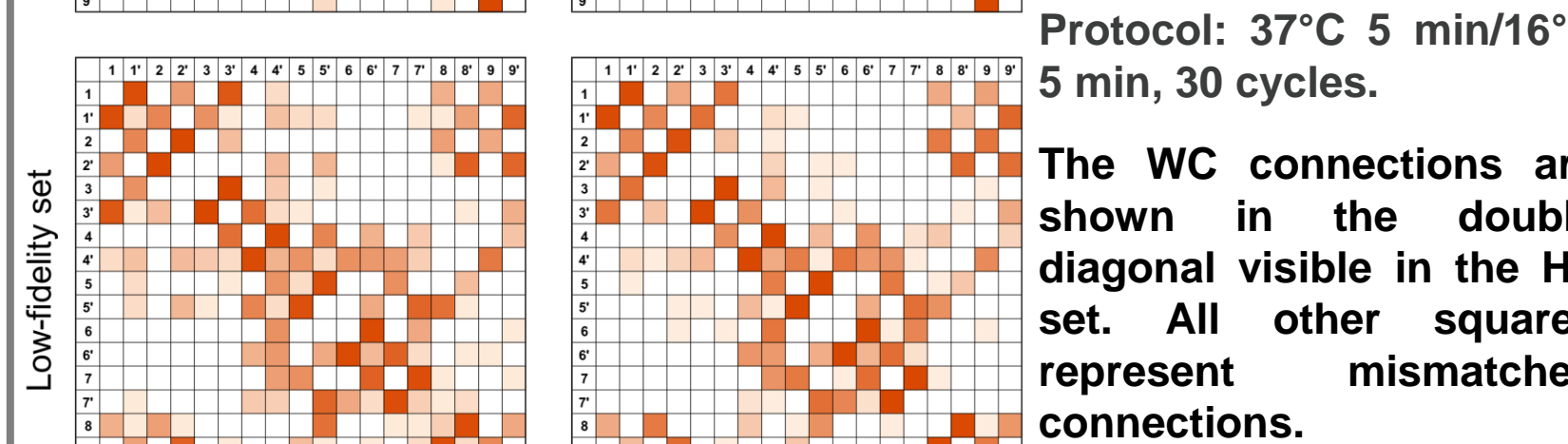
- Method allows comprehensive profiling of fidelity for a given ligase/end-structure pairing in one experiment.
- Overhangs range in fidelity from <50% WC to >99%.
- Most ligate with similar efficiency, but TNNA overhangs are notably inefficient.
- For even very low-fidelity overhangs, mismatch ligation products result from only a few (2 – 8) other ligation partners.
- Increasing temperature increases fidelity, but also increases bias against AT-rich overhangs.
- PEG drastically reduces fidelity.

Ligase Fidelity is Predictive in GGA

Overview of Golden Gate assembly design. Ten fragments of arbitrary sequence were designed, with junctions chosen to provide predicted high-fidelity or low-fidelity per the below table. A deletion prone set altered one HF junction to promote deletion/duplication of Fragment G, and a failure prone set including a predicted low efficiency junction. All products, full length and incomplete, were analyzed by sequencing to determine the identity of every junction in the assembly.

Junction	High-fidelity set	Low-fidelity set
1	AGGC	CGCC
2	ATCC	CGGG
3	ATCC	GCCA
4	TCAG	CGGT
5	AGCA	ACCC
6	TCCT	FGGG
7	AGTC	AGCC
8	TCAG	FGGG
9	ATCA	CGCC
10	TAGT	CGGG
11	CGCC	AGCA
12	CGCC	TCCT
13	CTGA	AGCC
14	GACT	FGCC
15	CGCA	CGCC
16	CGCT	CGCC
17	GGAA	AGCC
18	GCTT	TCGG

Protocol: 37°C 5 min/16°C 5 min, 30 cycles.

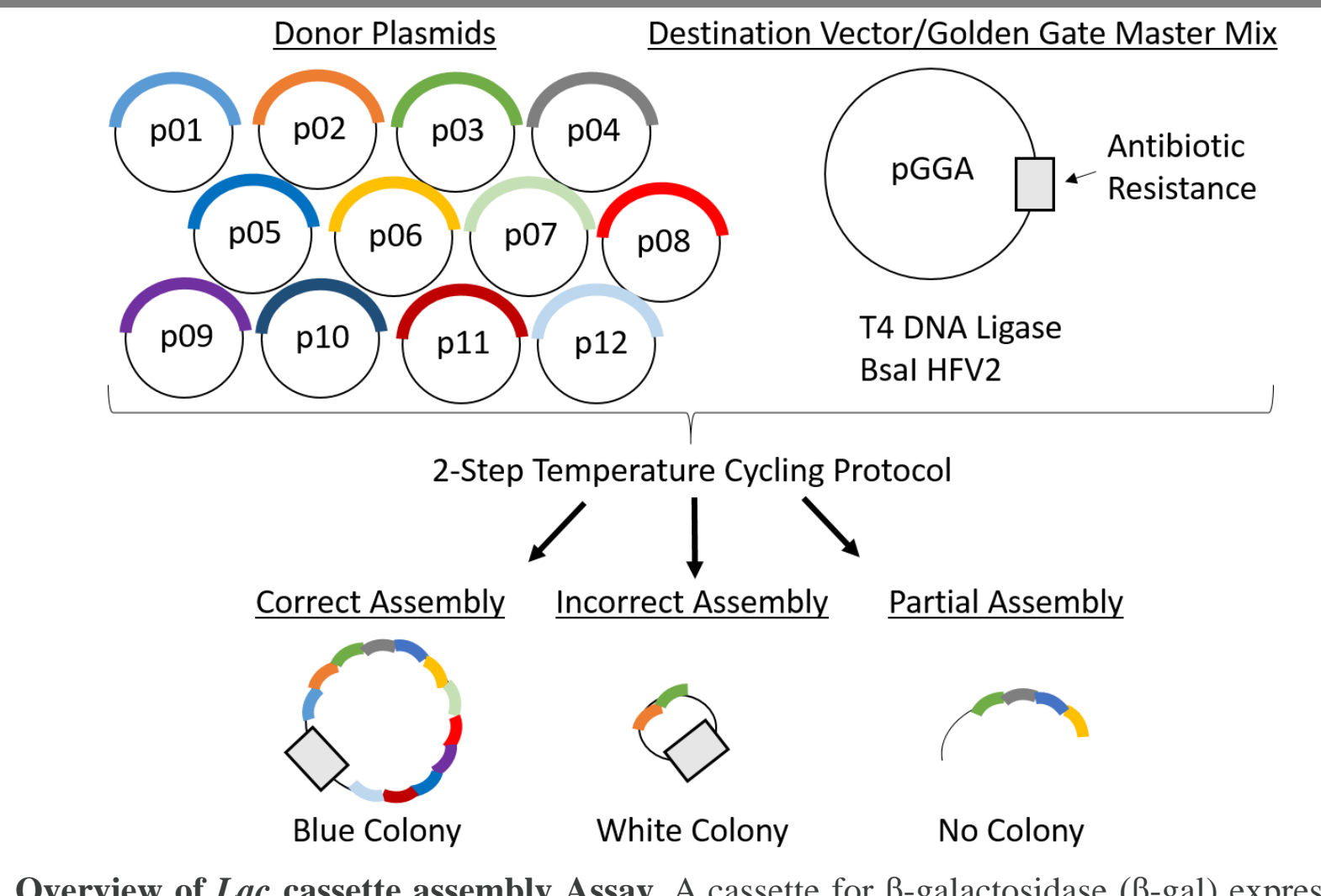


The WC connections are shown in the double diagonal visible in the HF set. All other squares represent mismatched connections. Predictions match results closely, excepting only lowest frequency events. Deletion prone set shows predicted deletion and duplications (indicated by arrows).

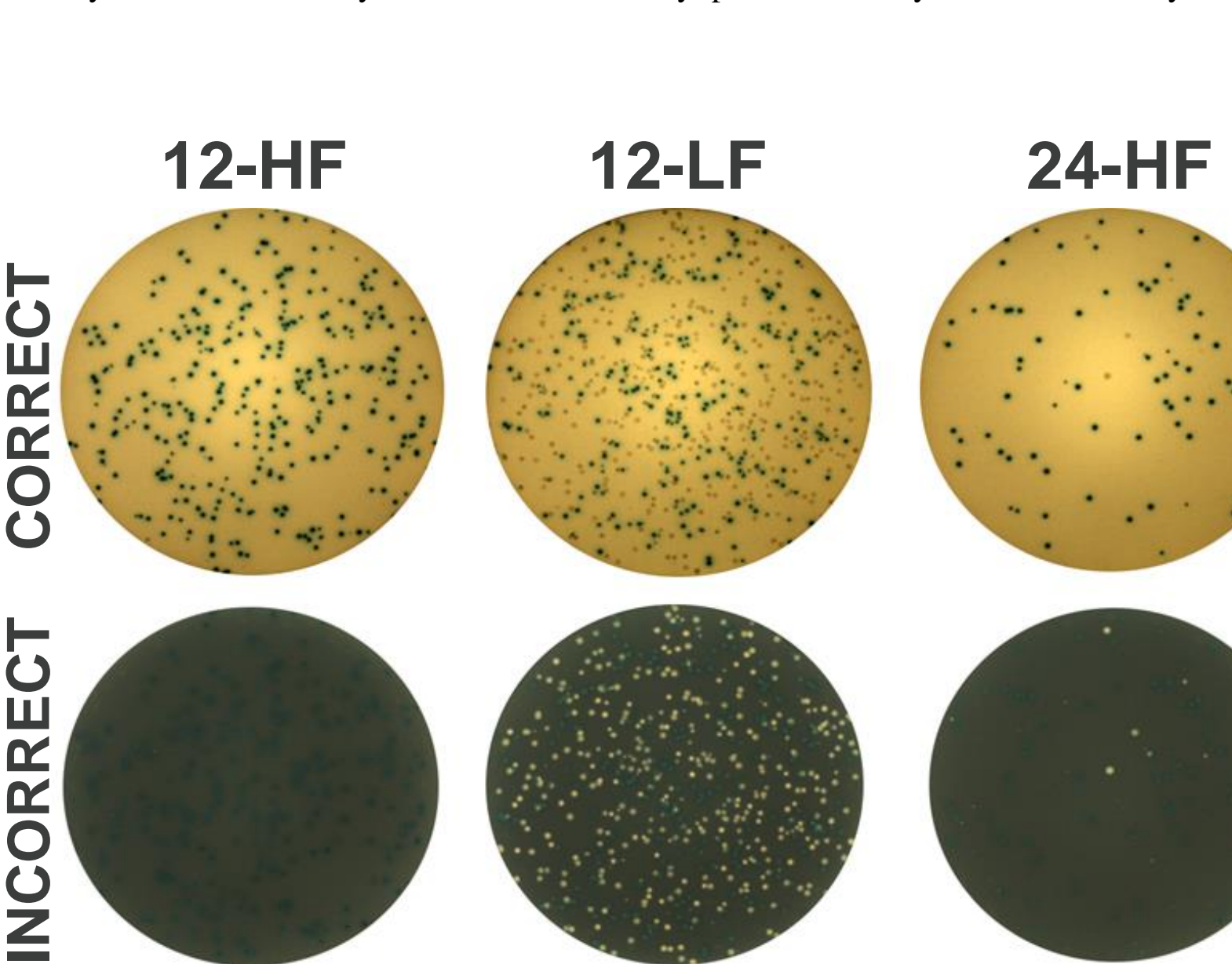
Failure prone set shows only a modest drop in connections at the low efficiency junction, less than predicted (indicated by arrows).

Bias against 100% GC overhangs observed. This result was not predicted by ligation assay, and is likely due to ligation profiling assay not accounting for cutting and melting steps.

12 and 24-Fragment One-Pot GGA



Overview of Lac cassette assembly Assay. A cassette for β-galactosidase (β-gal) expression was split into 12 or 24-fragments at junctions chosen based on the fidelity data and placed in holding vectors, with each fragment flanked by BsaI sites. In each assembly reaction, the fragments were combined with a pGGA destination vector containing the chloramphenicol resistance gene. A high-fidelity (HF) 12 and 24 fragment system was designed, as well as an intentionally deletion-prone low-fidelity (LF)12 fragment system. Assemblies containing all fragments in the proper order result in β-gal expression and a blue colony; incorrect assembly leads to a white colony; partial assembly leads to no colony.



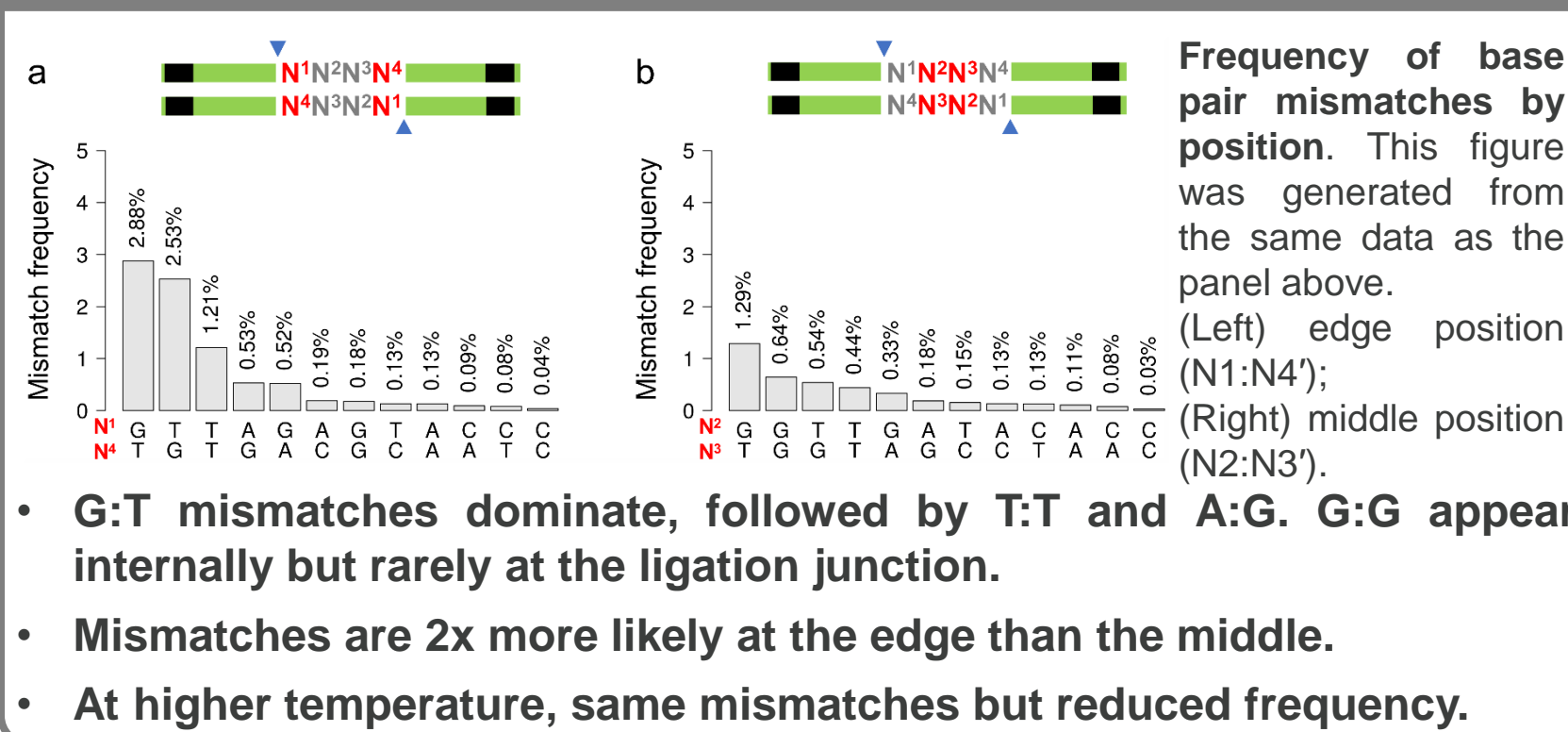
Protocol: 30 cycles (16°C 5 min, 37°C 5min), 5 min 60°C end hold. Results match predictions closely: 12-HF predicted 99% correct, observed 99.2 ± 0.6 % 12-LF predicted 31% blue colonies, observed 45 ± 5%. 24-HF predicted 91% blue colonies, observed 84 ± 5%.

Recommended Overhangs for GGA

Set	Number of overhangs	Estimated fidelity	Overhang sequences ¹
1 (MoClo Compatible)	15	98.5%	TGCC, GCAA, ACTA, TTAC, CAGA, TGTC, GAGC, AGGA, ATTC, CGAA, ATAG, AAGG, AACT, AAAA, ACCG
2	20	98.1%	AGTG, CAGG, ACTC, AAAA, AGAC, CGAA, ATAG, AACC, TACA, TAGA, ATGC, GATA, CTCG, GTAA, CTGA, ACAA, AGGA, ATTA, ACCG, CGCA
3	25	95.8%	CCTC, CTA, GACA, GCAC, AATC, GTAA, TGAA, ATTA, CCAG, AGGA, ACAA, TAGA, CGGA, GATA, CAGC, AACG, AAGT, CTCG, AGAT, ACCA, AGTG, GGTG, CGCA, AAAA, ATGA
4	30	91.7%	TACA, CTA, GGAA, GCCA, CACG, ACTC, CTTC, TCAA, GATA, ACTG, AACT, AAGC, CATA, GACC, AGGA, ATCG, AGAG, ATTA, CGGA, TAGA, AGCA, TGAA, ACAT, CCAG, GTGA, ACGA, ATAC, AAAA, AAGG, CAAC

¹ Only one member of each complementary overhang pair is shown. Set 1 is an extended MoClo set (TGCC, GCAA, ACTA, TTAC, CAGA, TGTC, GAGC) with additional 8 overhangs. All sets are predicted to assemble with the specified overall fidelity if every overhang and its complement is used; subsets of these sets are predicted to have even higher fidelity.

Position-Dependence of Mismatches



Frequency of base pair mismatches by position. This figure was generated from the same data as the panel above. (Left) edge position (N1:N4); (Right) middle position (N2:N3).

- G:T mismatches dominate, followed by T:T and A:G. G:G appear internally but rarely at the ligation junction.
- Mismatches are 2x more likely at the edge than the middle.
- At higher temperature, same mismatches but reduced frequency.

Summary/Conclusions

- Our new sequencing assay permits the rapid profiling of fidelity and bias in end-joining. Many sequences can be tested in a single pot.
- Fidelity and bias data for four-base overhangs predicts GGA accuracy.
- Specific erroneous junctions can be predicted in kind and degree
- Current data underestimates bias against 100% GC WC pairings.
- Use of data to guide junction overhang choice permits assembly of up to 24 fragments in one pot with high accuracy and efficiency.
- We have provided predicted best junction sets for GGA taking into account fidelity and efficiency.
- Future work will refine the data sets and examine bias of RE cutting.