



OligArch: A software tool to allow artificially expanded genetic information systems (AEGIS) to guide the autonomous self-assembly of long DNA constructs from multiple DNA single strands

Kevin M. Bradley^{1,2} and Steven A. Benner^{*1,2,3}

Full Research Paper

Open Access

Address:

¹Foundation for Applied Molecular Evolution, P.O. Box 13174, Gainesville FL 32604, USA, ²The Westheimer Institute for Science and Technology, 720 S. W. 2nd Avenue, Suites 201-208, Gainesville FL 32601, USA and ³Firebird Biomolecular Sciences LLC, 13709 Progress Blvd. Box 17, Alachua, FL 32615, USA

Email:

Steven A. Benner^{*} - sbenner@ffame.org

^{*} Corresponding author

Keywords:

AEGIS; bioinformatics; DNA self-assembly; long DNA constructs; software

Beilstein J. Org. Chem. **2014**, *10*, 1826–1833.

doi:10.3762/bjoc.10.192

Received: 02 February 2014

Accepted: 22 July 2014

Published: 11 August 2014

Editor-in-Chief: P. H. Seeberger

© 2014 Bradley and Benner; licensee Beilstein-Institut.

License and terms: see end of document.

Abstract

Synthetic biologists wishing to self-assemble large DNA (L-DNA) constructs from small DNA fragments made by automated synthesis need fragments that hybridize predictably. Such predictability is difficult to obtain with nucleotides built from just the four standard nucleotides. Natural DNA's peculiar combination of strong and weak G:C and A:T pairs, the context-dependence of the strengths of those pairs, unimolecular strand folding that competes with desired interstrand hybridization, and non-Watson–Crick interactions available to standard DNA, all contribute to this unpredictability. In principle, adding extra nucleotides to the genetic alphabet can improve the predictability and reliability of autonomous DNA self-assembly, simply by increasing the information density of oligonucleotide sequences. These extra nucleotides are now available as parts of artificially expanded genetic information systems (AEGIS), and tools are now available to generate entirely standard DNA from AEGIS DNA during PCR amplification. Here, we describe the OligArch (for "oligonucleotide architecting") software, an application that permits synthetic biologists to engineer optimally self-assembling DNA constructs from both six- and eight-letter AEGIS alphabets. This software has been used to design oligonucleotides that self-assemble to form complete genes from 20 or more single-stranded synthetic oligonucleotides. OligArch is therefore a key element of a scalable and integrated infrastructure for the rapid and designed engineering of biology.

Introduction

Automated synthesis of single stranded DNA fragments has, perhaps more than any other technology, enabled the development of "synthetic biology" as a modern field over the past 30 years [1-5]. While oligonucleotides can be reliably prepared by automated synthesis up to ca. 100 nucleotides in length and (even today) are most often used as primers, many seek to create large DNA (L-DNA) constructs by assembly of these fragments. Such engineered L-DNA might encode new and useful functions, including the manufacturing of biofuels, the synthesis of pharmaceuticals, and the development of new materials.

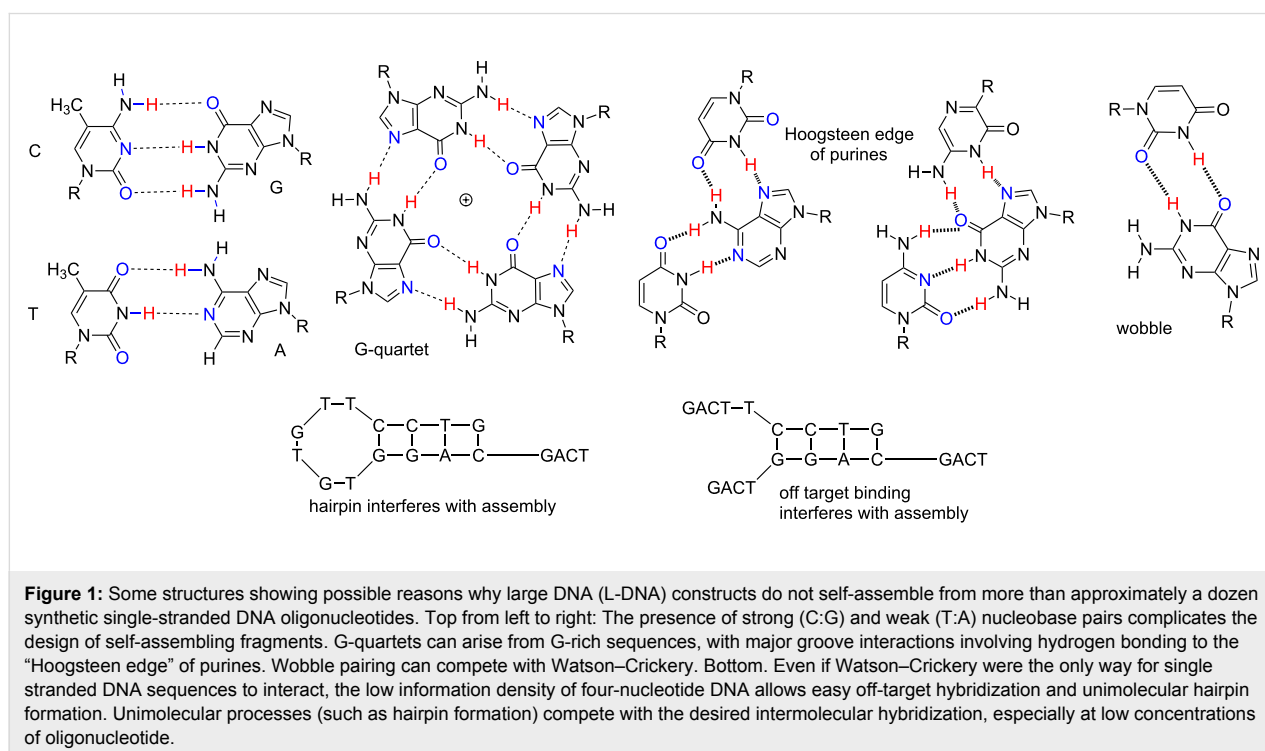
As it is taught to non-chemists, DNA appears to be an ideal molecule for such "hands-off" self-assembly. In this idealized "cartoon", two strands of DNA bind to each other perfectly, so long as their sequences are arranged so that A pairs with T and G pairs with C.

Even funding agencies have been captivated by this vision. For example, in 2011 the Army Research Office issued a small business grant solicitation seeking companies to design software to design 30,000 base pairs of single stranded DNA that would autonomously self-assemble to form nanostructures. In 2012, DARPA issued a small business grant solicitation seeking technology to assemble single-stranded synthetic fragments to give 20,000 base pair DNA constructs, essentially under this simple model for DNA behavior. More recently, DARPA has

initiated its "Foundries 1000" program, where large DNA "chassis" are hoped to occur in an entirely automated process.

If DNA in fact behaved according to this ideal, then the specificity of Watson–Crick nucleobase pairing might indeed allow autonomous self-assembly of an unlimited number of DNA strands to give L-DNA constructs of indefinitely large lengths. All that would be necessary is to design the requisite number of synthetic single strand fragments to remove off-target annealing, make them, and mix them. Once mixed, the designed strands would, in this view, simply fall together to form the target L-DNA structure.

Unfortunately, this simple model is also simplistic. With just four nucleotides, the information density of standard DNA is too low to allow (without exquisite design) even a dozen single strands to reliably self-assemble upon simple mixing. Further, even if rule-based Watson–Crick pairing were to be the only possible interaction, the combination of "strong" and "weak" G:C and A:T pairs makes design challenging. Also able to defeat self-assembly, DNA molecules can easily fold to give single-strand structures (such as hairpins), this unimolecular process competes with intermolecular duplex formation. Finally, a rich repertoire of non-Watson–Crick interactions (e.g., wobble, major groove binding) can compete with Watson–Crickery (Figure 1) to render autonomous self-assembly impossible.



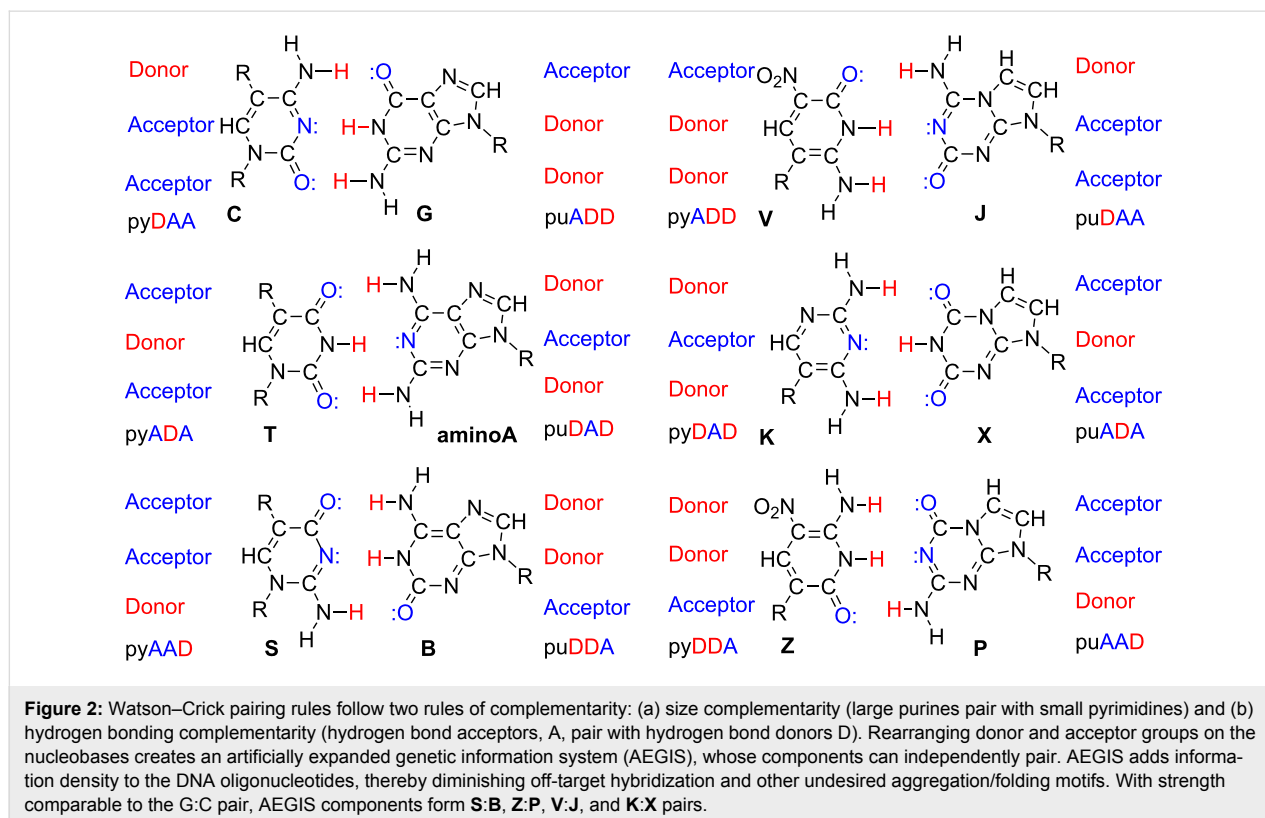
Even with great advances in recent years in large-scale DNA assembly using natural bases, the chance of failure of these assemblies grows very high once the number of fragments increases. With our own design that used fragments whose sequences were optimized for self-assembly [6], autonomous assembly typically failed between 16 and 24 oligonucleotides. Other methods of assembly, such as the Gibson Assembly [7] or SLIC [8], either limit the number of fragments to be used or rely on stepwise assembly of such syntheses, with the recommendation from Gibson that autonomous self-assembly of single stranded DNA fragments “be limited to perhaps a dozen fragments at a time”.

Fortunately, another development of synthetic biology offers an approach to mitigate these limitations of natural DNA as a matrix for autonomous self-assembly. This exploits a “second-generation” version of an artificially expanded genetic information system (AEGIS) [3,9]. AEGIS adds nucleotide building blocks to the four found in standard DNA (G, A, C, and T) by shuffling hydrogen-bonding units on the nucleobases, all while retaining the overall Watson–Crick nucleobase pairing geometry (Figure 2). These extra nucleotides bind to form additional nucleobase pairs orthogonally to the A:T and G:C pairs.

In principle, adding extra nucleotides in the genetic alphabet can mitigate the hybridization problems in highly complex

mixtures of single-strand DNA fragments simply by increasing the information density of the resulting DNA sequences. With four nucleotides, the number of possible 15mers (which form duplexes with convenient melting temperatures) is approximately 1.1 billion ($\approx 4^{15}$). While this number might appear to be large, it includes an enormous range of melting temperatures, a range arising because of the relative strengths of the G:C and A:T pair. Adding two additional nucleotides increases the number of potential hybridizing 15mers to 470 billion ($\approx 6^{15}$), nearly 500 fold higher. Adding four AEGIS nucleotides increases this number to 35 trillion ($\approx 8^{15}$). With a full AEGIS alphabet containing 12 nucleotide letters, approximately 1.5 quadrillion 15mers ($1.54 \times 10^{16} \approx 12^{15}$) are conceivable to serve as orthogonal hybridizing units. Further, AEGIS pairs are joined by three hydrogen bonds, giving them the strength of C:G pairs [3]. By increasing the information density of DNA, AEGIS DNA should more easily support “no hands” self-assembly by greatly increasing the number of possible unique fragments and maximize uniqueness of fragment ends.

Of course, a large DNA construct built with AEGIS nucleobases would, after it is assembled, still contain AEGIS components. For those who want an entirely natural end product, this is undesirable. Therefore, tools are needed to replace AEGIS pairs by standard pairs in processes that are rigorously rule-based.



Conversion of AEGIS nucleotides to standard nucleotides, it turns out, is facile by four AEGIS components: 2'-deoxy-5-methylisocytidine (trivially named **S**), 2'-deoxyisoguanosine (trivially named **B**), 2-amino-8-(1'- β -D-2'-deoxyribofuranosyl)-imidazo[1,2-*a*]-1,3,5-triazin-4(8*H*)one (trivially named **P**), and 6-amino-5-nitro-3-(1'- β -D-2'-deoxyribofuranosyl)-2(1*H*)-pyridone (trivially named **Z**) (Figure 2). In both cases, conversion is facilitated by forcing polymerases to mismatch AEGIS nucleotides with standard nucleotides by depriving the polymerase of the complementary AEGIS triphosphate. The specificity of mismatching is driven by intrinsic features of the AEGIS nucleobase. Thus, the **B**:T mismatch is enabled by a minor tautomeric form of **B**. The **Z**:G mismatch is enabled by the deprotonation of **Z**. The **P**:C mismatch is enabled by the protonation of **P** (Figure 3). These mismatches can occur both in vitro using PCR [10] and in vivo, in an engineered strand of *E. coli*. While rules for conversion have complexities, in their simplest forms, the **S**:**B** and **Z**:**P** pairs are converted to C:G and T:A pairs, respectively.

The availability of this new concept for the assembly of large DNA constructs from multiple inexpensive single-stranded oligonucleotides, together with the chemistry and enzymology needed to convert unnatural assemblies into entirely natural assemblies, creates the need for a software product to assist in the design of the fragments to be assembled. That product, OligArch (for “oligonucleotide architecting”) is described here. While other packages exist that output designed oligonu-

cleotides for large-scale synthesis, such as GeneGenie [11], DNAWorks [12], and Gene2Oligo [13], these all work with only natural bases, and require either a limited number of DNA fragments or a multistep process. OligArch is unique in allowing the use of AEGIS bases to greatly expand the number of fragments that can be used in a single-step self-assembly.

Discussion

Brief summary of OligArch software package

As input, the OligArch software package takes a sequence for a desired target DNA construct. As output, OligArch delivers sequences for a set of oligonucleotide fragments that include components of an artificially expanded genetic information system (AEGIS). OligArch designs these fragments so that the target DNA is produced after the fragments are annealed, after the annealed fragments are (optionally) extended by a DNA polymerase to fill in any gaps to give nicked DNA, after any nicks are sealed, and after the AEGIS pairs are replaced by standard pairs. OligArch also ensures that the increased information density of AEGIS oligonucleotides is exploited to avoid hairpin formation, off-target annealing, and undesired non-canonical structures. Thus, OligArch is a tool critical for exploiting the extra information density in AEGIS alphabets to assemble large DNA molecules.

The user of OligArch can designate within a given target DNA construct specific regions that encode proteins; OligArch ensures that the expressed protein is unchanged by the reassem-

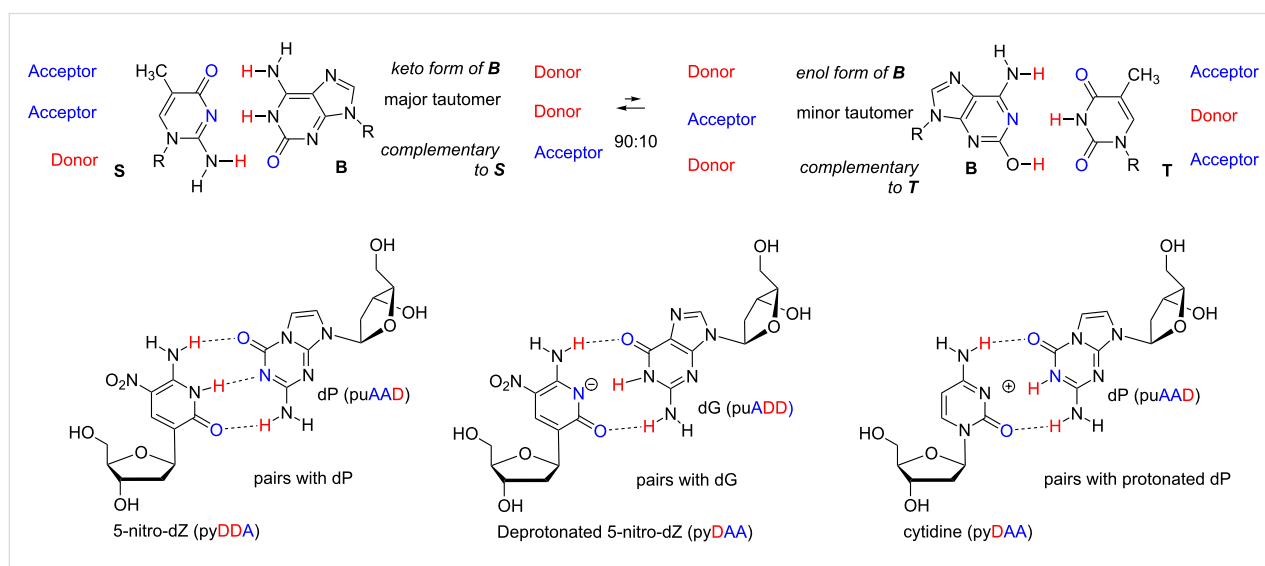


Figure 3: (top) The conversion of **S**:**B** pairs to **T**:**A** pairs involves tautomerization of **B** to give its minor enol form, which present a hydrogen bond Donor–Acceptor–Donor pattern complementary to **T**. If a strand containing **B** is copied by a polymerase that is not given any dSTP, mismatching of **T** opposite a minor enol tautomer of **B** leads (after two cycles of copying) to the replacements of **S**:**B** pairs by **T**:**A** pairs. (bottom) The conversion of **Z**:**P** pairs to **C**:**G** pairs involves the mismatching of **C** opposite a protonated **P**, and/or the mismatch of deprotonated **Z** opposite **G**. Thus, if a strand containing **Z** is copied at high pH by a polymerase that is not given any dPTP, mismatching of **G** opposite deprotonated **Z** leads (after two cycles of copying) to the replacements of **Z**:**P** pairs by **C**:**G** pairs. Conversely, if a strand containing **P** is copied at low pH by a polymerase that is not given any dZTP, mismatching of **G** opposite protonated **P** leads (after two cycles of copying) to the replacements of **Z**:**P** pairs by **C**:**G** pairs.

bled sequence. Further, the user can enter “sequence-dependent” regions (such as promoters, replication origins, etc.), and the application will ensure that the nucleotide sequence is completely unchanged in the reassembled sequence. Finally, short “non-changeable” regions (such as restriction enzyme recognition sites) can be entered to ensure no AEGIS substitutions occur at these locations.

User interface

Using OligArch’s web-based form, the user enters the sequence to be designed and specifies if the program should create oligonucleotides with short overlapping AEGIS spans (requiring extension by polymerase) or oligonucleotides with fully overlapping sequences. Also entered is the location of protein-coding, sequence-dependent, and non-changeable regions. The user then chooses design criteria such as optimal oligonucleotide length, longest and shortest acceptable oligonucleotides, number of AEGIS bases to use within AEGIS spans, AEGIS span melting temperatures (T_m), and other parameters.

Also customizable are criteria used for T_m calculations, including concentrations of oligonucleotides, Na^+ , and Mg^{2+} . The user may choose to use the AEGIS pairs **Z:P**, **S:B**, or **Z:P/S:B** for substitutions. Finally, the user can choose whether or not the designed sequence is circularized and whether or not oligos should be divided into individual assembly sets to allow step-wise assembly (only required for very large targets). Once the above information is submitted, OligArch runs using the algorithm described below. The designed oligonucleotides, along with any necessary warnings, are then displayed to the user in a table. Further, a detailed file can be downloaded with the designed oligonucleotides aligned to the original sequence in both tabular and graphical format.

Algorithm

OligArch attempts to create easily synthesizable oligos utilizing AEGIS base technology to ensure unique recognition sites exist in overlapping oligos. Once a user has input the sequence to be synthesized, sequence regions, and other design criteria, OligArch scans the sequence looking for positions where natural bases (A, G, T or C) could be substituted with AEGIS bases. These potential substitutions are stored along with the original bases in an indexed array. In designating potential substitutions, each of three categories of sequence regions have distinct rules:

1. Protein-coding regions: AEGIS nucleobases can be substituted only at the third site in codons where any nucleobase in the third position produces the same amino acid (Leu, Val, Ser, Pro, Thr, Ala, Arg, and Gly), as well

as codons where C/T or A/G are interchangeable in the third position (His, Gln, Asn, Lys, Asp, and Glu). The following AEGIS substitutions are allowed: **P** for a G, **Z** for a C, **B** for an A, and **S** for a T.

2. Sequence-dependent regions: AEGIS nucleobase substitutions in these regions are only allowed where known rules allow conversion back to ACTG bases results in the exact same sequence. Two such rules exemplify this type of substitution for **Z:P** pairs, which are **PP** to GG (and conversely **ZZ** to CC) and **PTP** to GTG (and conversely **ZAZ** to CAC) (Shaw & Benner, unpublished). For **S:B** pairs, any **S:B** pairs that are flanked on either side by a natural base, N, have 100% conversion back to T:A respectively (NSN to NTN and NBN to NAN)
3. Non-critical regions (default): Regions that are not critical to the sequence allow AEGIS substitution at any base. As with protein-coding regions, the following AEGIS substitutions are allowed: **P** for a G, **Z** for a C, **B** for an A, and **S** for a T.

The fourth category of sequence region, non-changeable regions, does not allow an AEGIS base substitution and is limited to a maximum size of 12 nucleotides.

Once this sequence substitution array has been created, OligArch searches for oligonucleotides that meet the user-entered specifications. Each oligonucleotide is composed of “AEGIS spans” on both ends, with optional AGTC sequence between each span. The AEGIS span is a sequence of nucleobases that contains AEGIS substitutions and that uniquely overlaps with its complement span. AEGIS spans attempt to have the optimal number of AEGIS–nucleobase substitutions, and must contain the minimal number of substitutions.

For short overlaps, the default size of the AEGIS spans is 12 to 15 nts with a required T_m between 46 and 58 °C, and the minimum number of AEGIS substitutions is two per span (with optimal set at four). For complete overlaps, the AEGIS spans cover half of the designed oligonucleotides, with a required T_m of 70 °C or greater, and a minimum of four AEGIS substitutions (with optimal set at 6). T_m ’s are calculated using the nearest-neighbor method along with unified entropy and enthalpy values from SantaLucia, et al. [14]. AEGIS bases **B** and **P**, both purines able to make 3 hydrogen bonds, are treated as G within T_m calculations, while **S** and **Z**, both pyrimidines able to form 3 hydrogen bonds, are treated as C. The exact T_m formula is:

$$T_m = \frac{\Delta H \frac{\text{kcal}}{\text{°C} \cdot \text{mol}}}{\Delta S + R \ln([\text{primer}] / 2)} - 273.15 \text{ °C} - \% \text{mismatch}$$

where R is the universal gas constant (1.987 cal/°C·mol) and ΔH and ΔS are the enthalpy and entropy of the base stacking adjusted for initiation factors [14], with ΔS also adjusted for salt concentrations using the formula [15]:

$$\Delta S[\text{Na}^+] = \Delta S[1 \text{ mol/L Na}^+] + 0.368 \cdot (\text{OligoLength} - 1) \cdot \ln\left(\left[\text{Na}^+\right] + \left[\text{Mg}^{++}\right] \cdot 140\right)$$

AEGIS segments are designed using a bidirectional sliding-window approach. Starting at the first position in the sequence, a sequence substitution array is used to design all possible AEGIS spans within the user-specified size range that match the input design criteria. AEGIS spans are then examined to include only those with an AEGIS nucleobase located within a user-specified number of bases from both the 5'- and 3'-end (the default is four nucleotides) to minimize unwanted complementarity. Remaining spans are then ranked based upon number of AEGIS nucleobases, distribution of AEGIS nucleobases, and presence and size of repeating nucleotide patterns.

Once ranked, these spans are compared against all other AEGIS spans being used in the current assembly set, with those that might hybridize with an unintended target being excluded. Potential hybridization is defined as any complementarity with a T_m of $\geq 50\%$ of the minimum span T_m that can be designed. Further, spans that will result in hairpin formation within designed primers are excluded. Hairpin formation is checked using the algorithm found in Primer3 [16,17] with a required max self-complementary score of 4 and additionally requiring a T_m of at least 6 °C in the stem sequence. Span comparisons continue until a single span is found that meets the above criteria. If no span can be created within the optimal window, the window is moved until an acceptable span can be created. With the bidirectional approach, the window moves 1 nt upstream from the optimal position, followed by 1 nt downstream, 2 nt upstream, 2 nt downstream, etc. until a valid AEGIS span can be found.

OligArch then attempts to design the next AEGIS span, which will be located downstream of the first span, in such a location that the two AEGIS spans (plus natural base sequence between the spans, if using short overlaps) creates an oligonucleotide of optimal length. The two AEGIS spans and any sequence between the spans become the first oligonucleotide, with the complement of the 2nd span becoming the start of the second nucleotide. This process continues until the entire sequence has been designed as individual nucleotides. On rare occasions, OligArch may need to design an oligonucleotide that is longer

or shorter than the user-specified criteria; a warning is issued if this occurs. For complete overlaps without gaps to fill, OligArch will also attempt to redesign the leading oligonucleotide, if necessary to allow proper design.

For circularized sequences, the complement of the first and last AEGIS span is used to create the final oligonucleotide sequences which will hybridize to circularize the product. If necessary, the program adjusts the lengths of the flanking fragments to ensure the final oligonucleotide is on the proper strand. For linear sequence, any sequence at the 5'- or 3'-end that could not be included in an AEGIS span is added to the end of the terminal oligonucleotides.

Results of OligArch designed assembly

The combination of the OligArch software with the ability to convert AEGIS nucleotides to standard nucleotides allows AEGIS to support the generation of large DNA constructs using the architecture described above. A representative assembly is shown in Figure 4, taken from the total synthesis of a gene encoding kanamycin resistance created using this strategy [6]. Here, oligonucleotides containing AEGIS components are designed by OligArch, with the protein coding region preserved to ensure the correct amino acid sequence is produced. These oligonucleotides are chemically synthesized, using chemistry entirely analogous to the chemistry used to synthesize standard DNA. The higher information density of the 6-letter DNA is then used to guide the autonomous assembly of large constructs by simple annealing. Gaps, if any, are filled in by polymerase extension, and the nicks in the product are sealed enzymatically with DNA ligase. After the higher information density provided by the AEGIS components has been exploited, the AEGIS components are removed via PCR to leave an entirely natural final DNA product, which was both complete and fully functional [6].

With the successful synthesis of the kanamycin resistance gene, we show that the AEGIS bases can be used for autonomous self-assembly of large-scale sequences, with conversion back to natural bases to preserve function of these sequences. With the theoretical ability to incorporate up to 8 AEGIS bases, we can far exceed the fragment count limits of natural base oligonucleotides inherent in current assembly methods [7,8], while also avoiding the need for step-wise assembly. With the implementation of OligArch, this technology will be easily harnessed for large scale sequence assembly.

Computer support and access to OligArch

OligArch is web-based application written in VB.net and utilizing ASP.net technology. It is hosted on a Dell R410 server running Windows Server 2008 R2 Standard. It is publicly

```

Forward Fragments:  CACCATGAGCCATATTCAACGGGAAACGTCGAGGCCGCGATTAAATTC AACATGGA$GCS$GAS$TT
Reverse Fragments:  CCTBCGBCTBAATATACCCATATTTA
AEGIS Construct: 1  CACCATGAGCCATATTCAACGGGAAACGTCGAGGCCGCGATTAAATTC AACATGGA$GCS$GAS$TTATATGGGTATAAAT 80
Converted Construct: 1  CACCATGAGCCATATTCAACGGGAAACGTCGAGGCCGCGATTAAATTC AACATGGATGCTGATTTATATGGGTATAAAT 80

Forward Fragments:  S GTCGGGCABT$B$GGTGC$GACAATCTATCGCTTGTATGGGAAGCCCGAS$GCGCC$B$GAG
Reverse Fragments:  CCCGAGCGCTATTBCAGCCCGT$S$AG$S
AEGIS Construct: 81  GGGCTCGCGATAAS$GTCGGGCABT$B$GGTGC$GACAATCTATCGCTTGTATGGGAAGCCCGAS$GCGCC$B$GAGTTGTTTCTG 160
Converted Construct: 81  GGGCTCGCGATAATGTCGGGC AATCAGGTGC$GACAATCTATCGCTTGTATGGGAAGCCCGATGCGCCAGAGTTGTTTCTG 160

Forward Fragments:  S GCCAAS$GAS$GT$S$ACAGATGAGATGGT$CAGACTAAACTGGCTGACGGAB$TTTATGCC$S$CT
Reverse Fragments:  TTTGTACCGTTTTCCATCGCAB$C$GGTT$B$CT$B$CAB
AEGIS Construct: 161  AAACATGGCAAAGGTAGCGT$S$GCCAAS$GAS$GT$S$ACAGATGAGATGGT$CAGACTAAACTGGCTGACGGAB$TTTATGCC$S$CT 240
Converted Construct: 161  AAACATGGCAAAGGTAGCGTTGCCAATGATGTTACAGATGAGATGGT$CAGACTAAACTGGCTGACGGAA$TTTATGCC$CT 240

Forward Fragments:  S C
Reverse Fragments:  S ACS$CS$GAS$GATGCATGGT$T$ACT$CACC$ACT$GCGAT$C$CC$G$B$AAA$C$B$G$C$B$T$T
AEGIS Construct: 241  SCCGACCATCAAGCATT$T$TATCCG$S$ACS$CS$GAS$GATGCATGGT$T$ACT$CACC$ACT$GCGAT$C$CC$G$B$AAA$C$B$G$C$B$T$T$C$C 320
Converted Construct: 241  TCCGACCATCAAGCATT$T$TATCCG$T$ACT$CCT$GAT$GAT$GCAT$GGT$T$ACT$CACC$ACT$GCGAT$C$CC$G$G$AAA$C$AG$C$ATT$C$C 320

Forward Fragments:  S GABAAS$ATT$G$S$GATGCGCTGGCAGT$GTT$CCTGCGCCGGT$T$GCA$S$T$CGAT$T
Reverse Fragments:  TCCATAATCTTCTTATAGGACTAAGTCC$B$CT$S$TT$B$TAAC$A$B$C
AEGIS Construct: 321  AGGTATTAGAA$AATATCCTGATTCAGG$S$GABAAS$ATT$G$S$GATGCGCTGGCAGT$GTT$CCTGCGCCGGT$T$GCA$S$T$CGAT$T 400
Converted Construct: 321  AGGTATTAGAA$AATATCCTGATTCAGGTGAAAATATTGTTGATGCGCTGGCAGT$GTT$CCTGCGCCGGT$T$GCA$T$T$CGAT$T 400

Forward Fragments:  C$T$G$T$S$T
Reverse Fragments:  TCG$S$CT$G$C$S$C$AGGC$G$CAAT$C$AC$GAAT$G$AAT$A$AC$G$S$T$T$G$G$T
AEGIS Construct: 401  CCTGT$S$T$G$T$A$A$T$G$T$C$T$T$T$A$A$C$A$G$C$G$A$T$C$G$C$G$T$A$T$T$T$C$S$T$C$G$C$S$C$AGGC$G$CAAT$C$AC$GAAT$G$AAT$A$AC$G$S$T$T$G$G$T 480
Converted Construct: 401  CCTGTTT$G$T$A$A$T$G$T$C$T$T$T$A$A$C$A$G$C$G$A$T$C$G$C$G$T$A$T$T$T$C$G$T$C$G$C$T$C$AGGC$G$CAAT$C$AC$GAAT$G$AAT$A$AC$G$G$T$T$G$G$T 480

Forward Fragments:  S$G$A$S$G
Reverse Fragments:  S$G$T$S$G$A$B$C$A$B$G$T$C$T$G$G$A$A$G$A$A$A$T$G$C$A$S$A$A$B$C$T$S$T$T$G$C$C$B$T
AEGIS Construct: 481  S$G$A$S$G$C$G$A$G$T$G$A$T$T$T$G$A$T$G$A$C$G$A$G$G$T$A$A$T$G$G$C$T$G$G$C$C$S$G$T$S$G$A$B$C$A$B$G$T$C$T$G$G$A$A$G$A$A$A$T$G$C$A$S$A$A$B$C$T$S$T$T$G$C$C$B$T 560
Converted Construct: 481  T$G$A$T$G$C$G$A$G$T$G$A$T$T$T$G$A$T$G$A$C$G$A$G$C$G$T$A$A$T$G$G$C$T$G$G$C$C$T$G$T$G$A$C$A$A$G$T$C$T$G$G$A$A$G$A$A$A$T$G$C$A$T$A$A$A$C$T$T$T$G$C$C$A$T 560

Forward Fragments:  A$S$T$T$C$T$B$C$T$S$G$A$S$A$A$C$C$T$A$T$T$T$T$G$A$C$G$A$G$G$G$A$A$T$A$A$T$A$G$G$T$T$G$T
Reverse Fragments:  A$G$A$G$T$G$G$C$T$A$A$G$T$C$A$G$C$A$G$T$G$A$G$T$A$C$C$A$C$T$B$A$A$G$A$G$S$G$A$B$C$T$B
AEGIS Construct: 561  T$C$T$C$A$C$C$G$G$A$T$T$C$A$G$T$C$G$T$C$A$C$T$C$A$T$G$G$T$G$A$S$T$T$C$T$B$C$T$S$G$A$S$A$A$C$C$T$A$T$T$T$T$G$A$C$G$A$G$G$G$A$A$T$A$A$T$A$G$G$T$T$G$T 640
Converted Construct: 561  T$C$T$C$A$C$C$G$G$A$T$T$C$A$G$T$C$G$T$C$A$C$T$C$A$T$G$G$T$G$A$T$T$T$C$A$C$T$T$G$A$A$C$C$T$A$T$T$T$T$G$A$C$G$A$G$G$G$A$A$T$A$A$T$A$G$G$T$T$G$T 640

Forward Fragments:  A$T$T$G$A$S$G$T$T$G$G$A$C$B$G$T$C$G$G$A$A$T$C$G$C$A$G$A$C$C$B$T$A$C$C$A$G$G$A$S$C$T$S$G$C$C$A$T$C$C$T$A$T$G$G$A$A$C$T$G$C$C$T$G$G$T$G$A$G$T$T$T$C$T$C$C
Reverse Fragments:  A$C$T$B$C$A$A$C$C$T$G$C$S$C$A$G$C$C$T$A$G$C$G$T$C$T$G$G$C$S$A$T$G$G$T$C$C$T$B$G$A$B$C
AEGIS Construct: 641  A$T$T$G$A$S$G$T$T$G$G$A$C$B$G$T$C$G$G$A$A$T$C$G$C$A$G$A$C$C$B$T$A$C$C$A$G$G$A$S$C$T$S$G$C$C$A$T$C$C$T$A$T$G$G$A$A$C$T$G$C$C$T$G$G$T$G$A$G$T$T$T$C$T$C$C 720
Converted Construct: 641  A$T$T$G$A$T$G$T$T$G$G$A$C$G$A$G$T$C$G$G$A$A$T$C$G$C$A$G$A$C$C$G$A$T$A$C$C$A$G$G$A$T$C$T$G$C$C$A$T$C$C$T$A$T$G$G$A$A$C$T$G$C$C$T$G$G$T$G$A$G$T$T$T$C$T$C$C 720

Forward Fragments:  S$T$C$B$T$T$A$C$A$G$A$A$B$C
Reverse Fragments:  C$S$G$A$S$A$T$G$A$A$S$A$A$B$T$T$G$C$A$G$T$T$C$A$T$T$G$A$T$G$C$T$C$G
AEGIS Construct: 721  S$T$C$B$T$T$A$C$A$G$A$A$B$C$G$G$C$T$T$T$T$C$A$A$A$A$A$T$A$T$G$G$T$A$T$G$A$A$A$T$C$C$S$G$A$S$A$T$G$A$A$S$A$A$B$T$T$G$C$A$G$T$T$C$A$T$T$G$A$T$G$C$T$C$G 800
Converted Construct: 721  T$T$C$A$T$T$A$C$A$G$A$A$A$C$G$G$C$T$T$T$T$C$A$A$A$A$A$T$A$T$G$G$T$A$T$G$A$A$A$T$C$T$G$A$T$A$A$A$A$A$T$T$G$C$A$G$T$T$C$A$T$T$G$A$T$G$C$T$C$G 800

Forward Fragments:  A$T$G$A$G$T$T$T$T$C$T$A$A$C$A$G$G$A$T$C$C$G$B$C$G$B$C$S$A$G
Reverse Fragments:  G$C$G$S$G$C$S$G$B$T$C$G$T$C$G$A$C$T$G$T$C$C$T$G$C$C$T$G$S$C$T$B$G$C$S$G$G$S$G$A$T$C
AEGIS Construct: 784  A$T$G$A$G$T$T$T$T$C$T$A$A$C$A$G$G$A$T$C$C$G$B$C$G$B$C$S$A$G$C$A$C$T$G$A$C$A$G$G$A$C$G$G$A$C$B$G$A$S$C$G$B$C$C$B$C$T$A$G 863
Converted Construct: 784  A$T$G$A$G$T$T$T$T$C$T$A$A$C$A$G$G$A$T$C$C$G$C$A$G$C$A$C$T$A$G$C$A$G$C$T$G$A$C$A$G$G$A$C$G$G$A$C$A$G$A$T$C$G$A$C$C$A$C$T$A$G 863

```

Figure 4: Representative assembly of oligonucleotides designed by OligArch built from the components of the six-nucleotide AEGIS GACTSB alphabet. The top two lines show the fragments. The line below shows the product before conversion. The bottom line shows the end product, entirely natural DNA, after conversion. This example, from [6], leads to the autonomous self-assembly of a gene that confers resistance to the antibiotic kanamycin.

accessible at <http://bioinformatics.ffame.org/bioinformatics/OligArch.aspx>.

Acknowledgements

We are indebted to the DARPA Foundries program (HR0011-12-C-0064) and, through Firebird Biomolecular Sciences LLC, the US Army (W911NF-12-C-0059) for partial support of this work. The basic research that made this work possible was funded by the Defense Threat Reduction Agency under its basic research program, including grant HDTRA1-13-1-0004.

References

- Ball, P. *Nature* **2004**, *431*, 624–626. doi:10.1038/431624a
- Benner, S. A.; Sismour, A. M. *Nat. Rev. Genet.* **2005**, *6*, 533–543. doi:10.1038/nrg1637
- Benner, S. A.; Yang, Z.; Chen, F. C. *R. Chim.* **2011**, *14*, 372–387. doi:10.1016/j.crci.2010.06.013
- Szybalski, W. *Control of Gene Expression*, Plenum Press: New York, 1974, pp. 23–24, 404–405, 411–412, 415–417.
- Rawis, R. L. *Chem. Eng. News* **2000**, *78*, 49–53. doi:10.1021/cen-v078n017.p049

6. Merritt, K. B.; Bradley, K. M.; Hutter, D.; Matsuura, M. F.; Rowold, D. R.; Benner, S. A. *Beilstein J. Org. Chem.* submitted.
7. Gibson, D. G.; Benders, G. A.; Andrews-Pfannkoch, C.; Denisova, E. A.; Baden-Tillson, H.; Zaveri, J.; Stockwell, T. B.; Brownley, A.; Thomas, D. W.; Algire, M. A.; Merryman, C.; Young, L.; Noskov, V. N.; Glass, J. I.; Venter, J. G.; Hutchison, C. A., III; Smith, H. O. *Science* **2008**, *319*, 1215–1220.
doi:10.1126/science.1151721
8. Li, M. Z.; Elledge, S. J. *Methods Mol. Biol. (N. Y., NY, U. S.)* **2012**, *852*, 51–59. doi:10.1007/978-1-61779-564-0_5
9. Benner, S. A. *Acc. Chem. Res.* **2004**, *37*, 784–797.
doi:10.1021/ar040004z
10. Yang, Z.; Chen, F.; Alvarado, J. B.; Benner, S. A. *J. Am. Chem. Soc.* **2011**, *133*, 15105–15112. doi:10.1021/ja204910n
11. Swainston, N.; Currin, A.; Day, P. J.; Kell, D. B. *Nucleic Acids Res.* **2014**, *42*, W395–W400. doi:10.1093/nar/gku336
12. Hoover, D. M.; Lubkowski, J. *Nucleic Acids Res.* **2002**, *30*, e43.
doi:10.1093/nar/30.10.e43
13. Rouillard, J.-M.; Lee, W.; Truan, G.; Gao, X.; Zhou, X.; Gulari, E. *Nucleic Acids Res.* **2004**, *32*, W176–W180. doi:10.1093/nar/gkh401
14. SantaLucia, J., Jr. *Proc. Natl. Acad. Sci. U. S. A.* **1998**, *95*, 1460–1465.
doi:10.1073/pnas.95.4.1460
15. von Ahsen, N.; Oellerich, M.; Armstrong, V. W.; Schütz, E. *Clin. Chem.* **1999**, *45*, 2094–2101.
16. Untergrasser, A.; Cutcutache, I.; Koressaar, T.; Ye, J.; Faircloth, B. C.; Remm, M.; Rozen, S. G. *Nucleic Acids Res.* **2012**, *40*, e115.
doi:10.1093/nar/gks596
17. Koressaar, T.; Remm, M. *Bioinformatics* **2007**, *23*, 1289–1291.
doi:10.1093/bioinformatics/btm091

License and Terms

This is an Open Access article under the terms of the Creative Commons Attribution License (<http://creativecommons.org/licenses/by/2.0>), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

The license is subject to the *Beilstein Journal of Organic Chemistry* terms and conditions: (<http://www.beilstein-journals.org/bjoc>)

The definitive version of this article is the electronic one which can be found at:
[doi:10.3762/bjoc.10.192](https://doi.org/10.3762/bjoc.10.192)