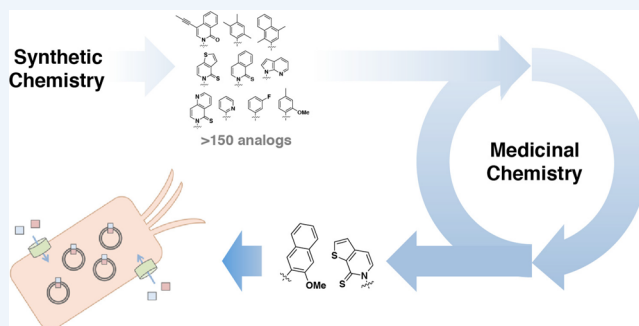


# Expansion of the Genetic Alphabet: A Chemist's Approach to Synthetic Biology

Aaron W. Feldman<sup>1</sup> and Floyd E. Romesberg<sup>1\*</sup>

Department of Chemistry, The Scripps Research Institute, La Jolla, California 92037, United States

**CONSPECTUS:** The information available to any organism is encoded in a four nucleotide, two base pair genetic code. Since its earliest days, the field of synthetic biology has endeavored to impart organisms with novel attributes and functions, and perhaps the most fundamental approach to this goal is the creation of a fifth and sixth nucleotide that pair to form a third, unnatural base pair (UBP) and thus allow for the storage and retrieval of increased information. Achieving this goal, by definition, requires synthetic chemistry to create unnatural nucleotides and a medicinal chemistry-like approach to guide their optimization. With this perspective, almost 20 years ago we began designing unnatural nucleotides with the ultimate goal of developing UBPs that function *in vivo*, and thus serve as the foundation of semi-synthetic organisms (SSOs) capable of storing and retrieving increased information. From the beginning, our efforts focused on the development of nucleotides that bear predominantly hydrophobic nucleobases and thus that pair not based on the complementary hydrogen bonds that are so prominent among the natural base pairs but rather via hydrophobic and packing interactions. It was envisioned that such a pairing mechanism would provide a basal level of selectivity against pairing with natural nucleotides, which we expected would be the greatest challenge; however, this choice mandated starting with analogs that have little or no homology to their natural counterparts and that, perhaps not surprisingly, performed poorly. Progress toward their optimization was driven by the construction of structure–activity relationships, initially from *in vitro* steady-state kinetic analysis, then later from pre-steady-state and PCR-based assays, and ultimately from performance *in vivo*, with the results augmented three times with screens that explored combinations of the unnatural nucleotides that were too numerous to fully characterize individually. The structure–activity relationship data identified multiple features required by the UBP, and perhaps most prominent among them was a substituent ortho to the glycosidic linkage that is capable of both hydrophobic packing and hydrogen bonding, and nucleobases that stably stack with flanking natural nucleobases *in lieu* of the potentially more stabilizing stacking interactions afforded by cross strand intercalation. Most importantly, after the examination of hundreds of unnatural nucleotides and thousands of candidate UBPs, the efforts ultimately resulted in the identification of a family of UBPs that are well recognized by DNA polymerases when incorporated into DNA and that have been used to create SSOs that store and retrieve increased information. In addition to achieving a longstanding goal of synthetic biology, the results have important implications for our understanding of both the molecules and forces that can underlie biological processes, so long considered the purview of molecules benefiting from eons of evolution, and highlight the promise of applying the approaches and methodologies of synthetic and medical chemistry in the pursuit of synthetic biology.



## INTRODUCTION

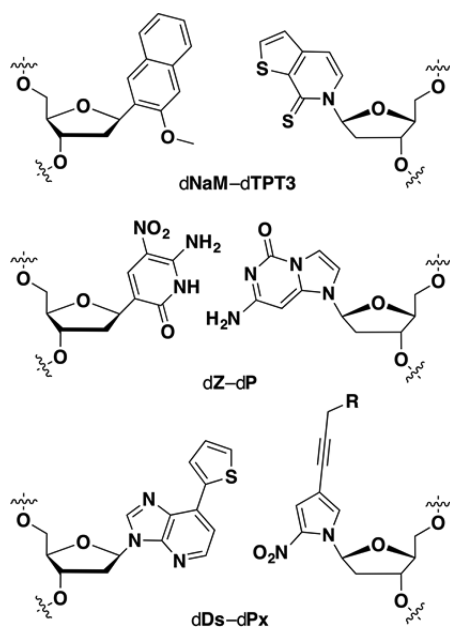
The field of synthetic biology was first defined in 1911 by Stéphane Leduc<sup>1</sup> with the goal of creating new biological forms and functions. The modern field is largely focused on using the engineering-like approach of design, test, and standardize to optimize naturally derived “parts”, most commonly novel DNA elements. However, the most fundamental approach to create new forms and functions is to develop unnatural base pairs (UBPs) that expand the genetic alphabet and, thus, increase the amount of information that may be stored in an organism’s DNA. The effort to expand the genetic alphabet was first championed by Steven Benner and focused on unnatural nucleotides bearing hydrogen-bonding (H-bonding) patterns that are orthogonal to those employed by the natural base pairs (Figure 1).<sup>2</sup> However, it was unclear if H-bonding was the only

force suitable for controlling base pairing. Indeed, the Kool laboratory reported the remarkable observation that a DNA polymerase could selectively insert the difluorotoluene analog of dTTP opposite dA in the template.<sup>3</sup> We and the Hirao laboratory thus initiated efforts to use hydrophobic and packing forces to control UBP formation. Hirao has focused on derivatizing natural purine and pyrimidine scaffolds to create “shape complementary” UBPs (Figure 1),<sup>4</sup> while we focused on nucleotides bearing nucleobase analogs with little to no homology to their natural counterparts.

Although our efforts were consistent with the tenets of modern synthetic biology, we used synthetic chemistry to

Received: August 16, 2017

Published: December 2, 2017



**Figure 1.** dNaM-dTPT3, dZ-dP, and dDs-dPx (R = H or  $-\text{CH}(\text{OH})-\text{CH}_2\text{OH}$ ) UBPs.

generate the required parts, which is in many ways more consistent with Leduc's original vision.<sup>1</sup> Furthermore, we optimized the parts using a medicinal chemistry-like approach, inspired by the perspective that any effort endeavoring to develop foreign, man-made parts that function within a cell will need to optimize solubility, cellular uptake, stability, off-target activity, toxicity, dose, and dosing regimen. The field of medicinal chemistry approaches the same problems through the construction of structure–activity relationships (SARs) that allow for empirical optimization in the absence of a complete understanding of the process being optimized, a strategy that we adapted for UBP development.

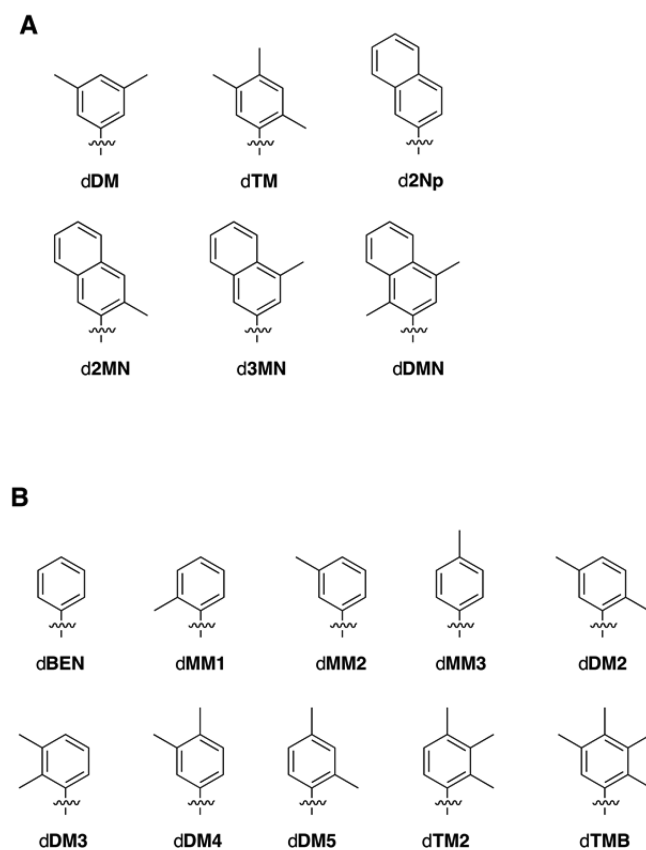
Here, we recount our efforts to develop UBPs that ultimately culminated in the discovery of the dNaM-dTPT3 UBP (Figure 1), as well as a family of related UBPs, all of which have now been used to create semi-synthetic organisms (SSOs) that store<sup>5,6</sup> and retrieve<sup>7</sup> increased information. This places us at the doorstep of realizing Leduc's vision of creating organisms with novel forms and functions.

### ■ PARTS OPTIMIZATION: FIRST GENERATION UBP CANDIDATES

The synthetic biology parts required for the expansion of the genetic alphabet are the unnatural nucleotides that selectively pair to form a UBP, and their optimization, at least initially, was measured by both duplex stability and DNA polymerase recognition. In most cases, analogs were synthesized as triphosphates (referred to as dXTP, where X is a nucleoside analog), as well as phosphoramidites for incorporation into DNA via solid phase synthesis. While characterization of the stability of duplex DNA containing the UBPs revealed many interesting trends, such as the effects of solvation,<sup>8</sup> stability proved uncorrelated with polymerase recognition, as is also the case with natural base pairs,<sup>9</sup> and polymerase recognition quickly emerged as our primary focus. The majority of our early efforts employed the exonuclease-deficient Klenow fragment of *E. coli* DNA polymerase I (Kf). We focused specifically on two steps, the insertion of an unnatural triphosphate opposite its

cognate analog in a template (a step that we also refer to as UBP synthesis or unnatural nucleotide incorporation), and the continued primer extension by insertion of the next correct triphosphate, in each case characterizing efficiency (second-order rate constant,  $k_{\text{cat}}/K_{\text{M}}$ ) and fidelity ( $(k_{\text{cat}}/K_{\text{M}})_{\text{correct}}/(k_{\text{cat}}/K_{\text{M}})_{\text{incorrect}}$ ). Steady-state conditions were employed, which allowed rapid SAR construction but only provided information about the rate-limiting step of insertion or extension, both of which are actually a complex series of reactions (e.g., triphosphate binding, conformational changes, phosphoryl transfer, product release). During early development, this was not a problem as the rate-limiting step was invariably phosphoryl transfer.

Our search began with the simple benzene and naphthalene nucleobase analogs (Figure 2A).<sup>10,11</sup> We found that dDMTP is

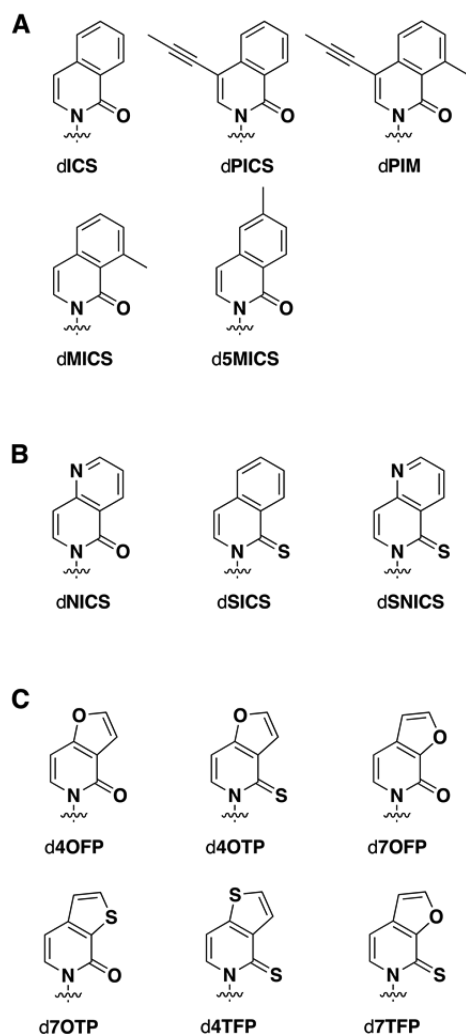


**Figure 2.** (A) First generation methylated benzene and naphthalene analogs. (B) Second generation methylated benzene analogs. Sugar and phosphate groups omitted for clarity.

a poor polymerase substrate, as incorporation opposite dDM or dTM in the template was virtually undetectable. In contrast, dTMTP was more efficiently incorporated opposite dDM ( $k_{\text{cat}}/K_{\text{M}} = 1.4 \times 10^6 \text{ M}^{-1} \text{ min}^{-1}$ ) and dTM ( $k_{\text{cat}}/K_{\text{M}} = 2.2 \times 10^6 \text{ M}^{-1} \text{ min}^{-1}$ ), only 20–30-fold less than the efficiency with which a natural base pair is synthesized in the same sequence context. However, both the dTM-dDM heteropair and dTM-dTM self-pair are limited by poor extension and mispairing with dA, likely because dA is the most hydrophobic of the natural nucleotides. We also found that d2Np efficiently directs the insertion of d2NpTP ( $k_{\text{cat}}/K_{\text{M}} = 2.8 \times 10^6 \text{ M}^{-1} \text{ min}^{-1}$ ), but misincorporation of dATP was again problematic ( $k_{\text{cat}}/K_{\text{M}} = 1.1 \times 10^5 \text{ M}^{-1} \text{ min}^{-1}$ ). Addition of a methyl substituent to the position ortho to the glycosidic bond (hereafter referred to

simply as the ortho position) yielded d2MN, which generally increased the rate of incorporation of the triphosphate against hydrophobic analogs in the template. However, d2MNTP is again efficiently inserted opposite dA. Moving the single methyl substituent from the ortho to the meta position (d3MN) reduces insertion opposite dA but also reduces the efficiency of UBP synthesis. The additional methyl substituent of dDMN resulted in efficient pairing with both dTM and d2MN, but it also increased the rate of dATP insertion when in the template. Thus, while mispairing generally remained problematic, at this point it was clear that several of these UBPs showed promising synthesis rates. However, none of them supported continued primer extension at a detectable level ( $k_{\text{cat}}/K_{\text{M}} < 10^3 \text{ M}^{-1} \text{ min}^{-1}$ ).

The isocarbostyryl scaffold also received early development attention (Figure 3A,B).<sup>10,12–14</sup> When in the template, the



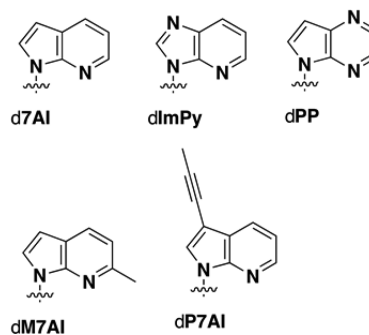
**Figure 3.** (A) Isocarbostyryl analogs. (B) Heteroatom derivatized isocarbostyryl analogs. (C) Furan and thiophene fused pyridone and thiopyridone analogs. Sugar and phosphate groups omitted for clarity.

simplest of the series, dICS, directs the insertion of triphosphates bearing simple substituted benzene nucleobases, such as dDMTP, dTMTP, or dDMNTP, with only modest efficiency, but d2MNTP is inserted significantly more efficiently. Addition of a 1-propynyl group to the 7-position of dICS, affording dPICS, results in a triphosphate that is

generally inserted more efficiently opposite other unnatural analogs. The methyl group of dPIM resulted in efficient but indiscriminate insertion of the triphosphate, while addition of the methyl group of dMICS resulted in the indiscriminate templating of natural triphosphates. Although the methyl group of d5MICS had little systematic effect on UBP synthesis, it did dramatically reduce self-pairing. Despite variable rates of UBP synthesis, the rates of extension remained problematic. While 6-aza substitution (dNICS) results in 2-fold reduced self-pair synthesis, it also results in a 2-fold increased rate of extension. Remarkably, replacing the oxygen of dICS with sulfur (dSICS) results in an 80-fold increase in the rate of self-pair synthesis and a 4-fold increase in the rate of extension. The combination of both modifications in dSNICS retained the increased rate of synthesis but further increased the rate of extension 12-fold ( $k_{\text{cat}}/K_{\text{M}} = 2.2 \times 10^4 \text{ M}^{-1} \text{ min}^{-1}$ ). These results provided early hints as to the important but complicated role of the ortho substituent.

A series of pyridones and thiopyridones fused to furan and thiophene rings in both meta- and para-linked orientations were investigated (Figure 3C).<sup>15</sup> UBP synthesis with these analogs was generally inefficient, with a para-linked furan appearing to be particularly detrimental. The thiopyridone triphosphate analogs are generally inserted more efficiently, but the effect of sulfur was less pronounced than with the isocarbostyryl scaffold, and none of these analogs emerged as particularly promising.

Nucleoside analogs bearing azaindole scaffolds (Figure 4) were found to efficiently pair with various unnatural nucleotides

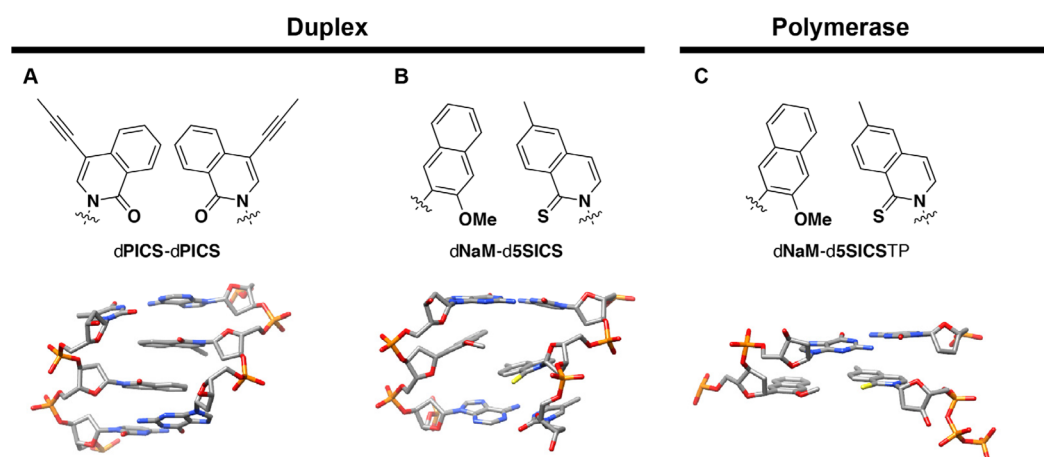


**Figure 4.** Azaindole analogs. Sugar and phosphate groups omitted for clarity.

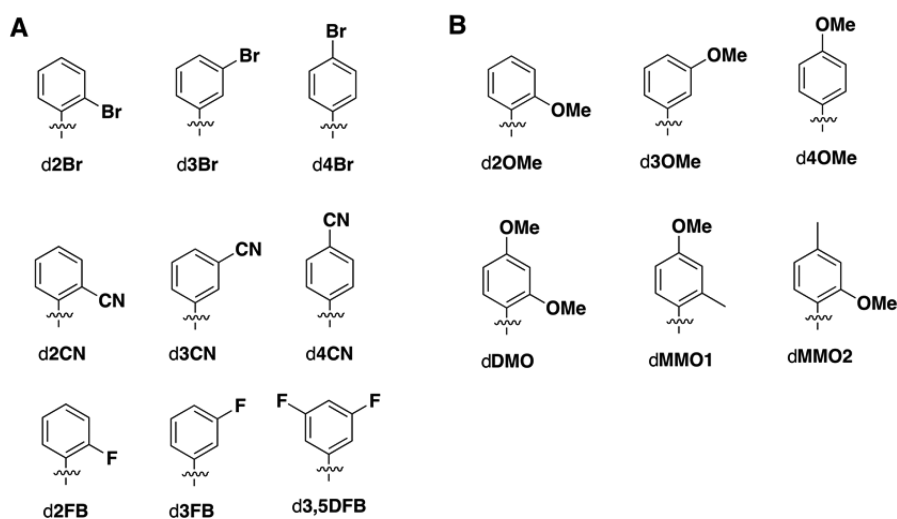
in the template, and with reasonable selectivity against the natural triphosphates when in the template, with the exception of dP7AI, which efficiently templates the insertion of dCTP.<sup>16</sup> However, while UBP synthesis again proved amenable to optimization, the resulting UBPs were generally refractory to extension, which by now had emerged as the greatest challenge to optimization.

## ■ A STRUCTURAL INTERLUDE

The dominant SAR that emerged from these first generation UBP candidates was that while the hydrophobicity and aromatic surface area of aromatic nucleobase derivatives appeared promising for the optimization of UBP synthesis, the resulting UBPs generally proved refractory to continued primer extension. The NMR structure of a duplex containing the dPICS self-pair (Figure 5A) was solved to better understand this SAR.<sup>17</sup> These studies revealed a generally unperturbed duplex structure, with the propargyl moieties of dPICS disposed as expected in the major groove; however,



**Figure 5.** (A) Duplex structure of DNA containing the dPICS–dPICS UBPs.<sup>16</sup> (B) Duplex structure of DNA containing the dNaM–d5SICS UBPs.<sup>36</sup> (C) Structure of d5SICSSTP paired opposite dNaM in the polymerase active site.<sup>37</sup> In chemical structures, sugar and phosphate groups omitted for clarity.



**Figure 6.** (A) Bromo-, cyano-, and fluoro-substituted benzene analogs. (B) Methoxy-substituted benzene analogs. Sugar and phosphate groups omitted for clarity.

considerable distortion was observed at the site of the UBPs itself. Rather than adopting the canonical edge-on Watson–Crick geometry, the dPICS nucleobases interact via cross-strand intercalation. This intercalation appears to be driven by favorable stacking of the large aromatic surfaces of each dPICS nucleobase and has been observed with other nucleotides bearing nucleobase analogs with extended aromatic surface area.<sup>18</sup> We hypothesized that the same mode of pairing occurs at the primer terminus, which we envisioned would account for the efficient rates of UBPs synthesis, but also the inefficient rates of continued extension, as deintercalation would be required to appropriately position the primer terminus.

### CONTINUED PARTS OPTIMIZATION: SECOND GENERATION UBPs CANDIDATES

Based on the SAR generated with the first generation analogs, we began testing whether smaller nucleobase analogs could be optimized for UBPs synthesis, with the expectation that they would be less prone to cross-strand intercalate. These second generation efforts started with a more complete analysis of the simple benzene scaffold explored previously (Figure 2B). The parent analog, dBENTP, is poorly recognized by Kf.<sup>19</sup> We then

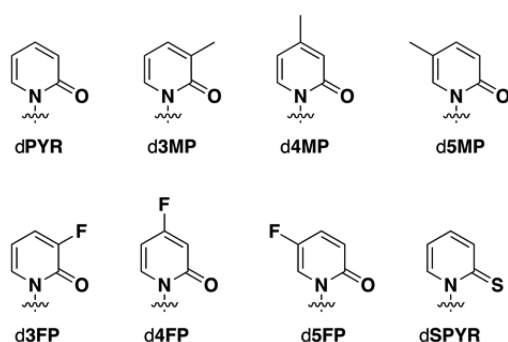
explored dMM1, dMM2, and dMM3, but found little improvement.<sup>19</sup> The efficiency of triphosphate insertion was progressively increased with dDMTP, dDM2TP, dDM3TP, dDM4TP, and dDM5TP, and insertion opposite dA was eliminated with dTMTP. Despite this progress with insertion, extension of primers terminating with these analogs remained inefficient.

We next explored heteroatom derivatization of these small scaffolds with bromo-, fluoro-, and cyano-substituents (Figure 6A). Bromo- and cyano-substituents tended to decrease the rate of mispairing with natural nucleotides,<sup>20</sup> but extension rates remained poor. A systematic analysis of fluoro substitution identified a single meta substituent (d3FB) as particularly interesting, with the resulting self-pair both synthesized and extended with at least moderate efficiency,<sup>21</sup> but further optimization efforts proved unproductive.

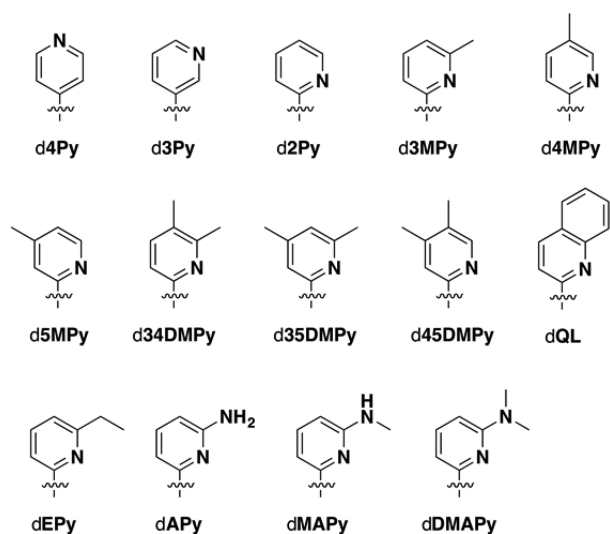
The situation improved with a family of methoxy-derivatized analogs (Figure 6B), and those possessing an ortho-methoxy substituent were the first to provide UBPs that were consistently extended with at least reasonable efficiency.<sup>22</sup> In fact, with dMMO2 paired against dTM in the template, the resulting primer was extended with an efficiency that is only 36-

fold lower than that of a natural base pair in the same sequence context. Mutation of the polymerase indicated that the increased extension resulted from the ability of the ortho-methoxy group to accept an H-bond, and similar substituents within the natural nucleobases are known to play a similar role.<sup>23–25</sup> Consistent with the need for an ortho H-bond acceptor at the primer terminus, dTM paired opposite dMMO2 in the template is only extended poorly ( $k_{\text{cat}}/K_{\text{M}} = 6.3 \times 10^3 \text{ M}^{-1} \text{ min}^{-1}$ ). Although the d3FB self-pair was an exception, the SAR strongly suggested that an ortho H-bond acceptor was essential, and its inclusion emerged as a central design theme.

A variety of heterocyclic N- and C-nucleotides, which can also position an H-bond acceptor at the position ortho to the glycosidic linkage, were next examined<sup>26–28</sup> (Figures 7 and 8).



**Figure 7.** Derivatized monocyclic pyridone analogs. Sugar and phosphate groups omitted for clarity.



**Figure 8.** Pyridine and substituted pyridine analogs. Sugar and phosphate groups omitted for clarity.

The pyridone analog triphosphates were inserted with only marginal efficiency, but the resulting UBPs were extended with greater efficiency, consistent with the proposed role of the ortho H-bond acceptor. Conversion of dPYR to the corresponding thiopyridone (dSPYR) resulted in UBPs that were still reasonably well extended but synthesized less efficiently.<sup>29</sup> No pyridine analogs were efficiently inserted as triphosphates, but d2Py was better extended when at the primer terminus than was d3Py or d4Py, again consistent with the importance of an H-bond acceptor in the developing minor groove.<sup>28</sup> Problematically, however, the pyridine analogs paired

more efficiently with dATP than any unnatural triphosphates when in the template, and no improvements were found with various alkyl or heteroatom substituents or with an increased aromatic surface (dQL).<sup>30</sup>

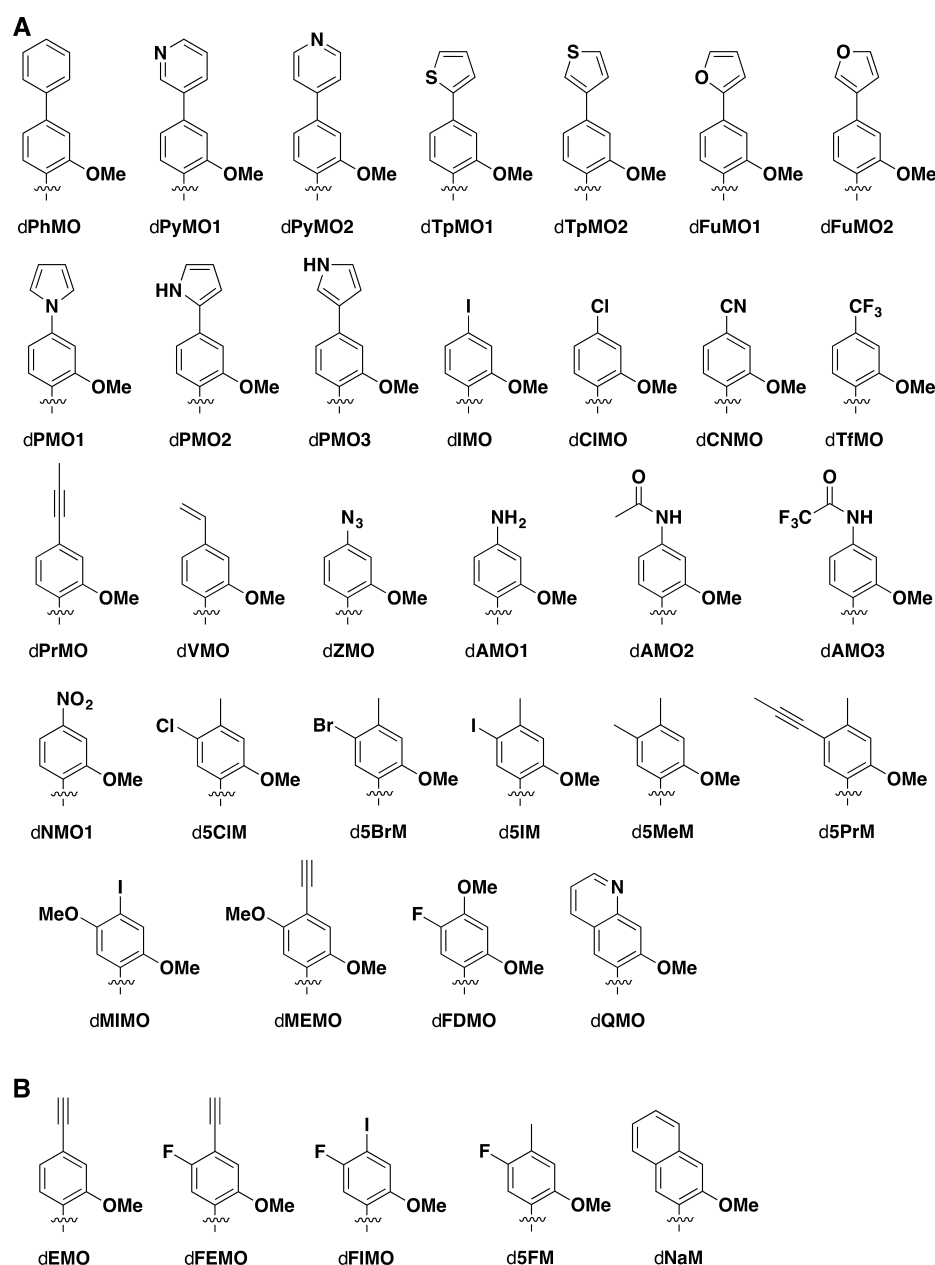
The most pronounced SAR generated from these second generation candidates was that efficient UBP synthesis requires a hydrophobic ortho substituent, while efficient extension requires the same substituent to be more hydrophilic, at least hydrophilic enough to act as an H-bond acceptor. This apparent physicochemical dichotomy challenged our confidence that both UBP synthesis and extension could be simultaneously optimized.

## ■ PARTS OPTIMIZATION: THIRD GENERATION UBP CANDIDATES

With no clear strategy to satisfy the conflicting demands of UBP synthesis and extension, we pivoted to a screen-based strategy. Two complementary screens were performed, one gel-based screen that focused on UBP incorporation and extension under steady state conditions, and one fluorescence-based plate screen that focused on the efficiency and fidelity of full length synthesis.<sup>29</sup> Remarkably, from 3600 candidate UBPs, both screens identified dMMO2–dSICS as the most promising. This UBP appeared to satisfy the physicochemical contradiction, because the sulfur of the dSICS thioamide is relatively hydrophobic but still able to accept an H-bond, while simple bond rotation allows the methoxy group of dMMO2 to direct either a hydrophobic methyl group or the oxygen lone pairs into the developing minor groove.

The identification of dMMO2–dSICS reinvigorated our design efforts. Steady-state kinetics revealed that while dMMO2 and dSICS were both incorporated and extended relatively efficiently ( $k_{\text{cat}}/K_{\text{M}} = 3.4 \times 10^5 \text{ M}^{-1} \text{ min}^{-1}$  and  $1.7 \times 10^6 \text{ M}^{-1} \text{ min}^{-1}$ , respectively, with dSICS in the template, and  $1.4 \times 10^6 \text{ M}^{-1} \text{ min}^{-1}$  and  $1.1 \times 10^6 \text{ M}^{-1} \text{ min}^{-1}$ , respectively, with dMMO2 in the template), replication fidelity was limited by dSICS self-pairing. Based on previous SAR, we explored the addition of a methyl group to the distal aromatic ring of dSICS to introduce steric interactions to disfavor self-pairing. These efforts identified d5SICS, and thus dMMO2–d5SICS emerged as our lead UBP.

We next turned to the optimization of dMMO2TP insertion opposite d5SICS, which was now the rate-limiting step of replication (Figure 9). Based on previous results that a meta-fluoro substituent or an expansion of aromatic surface area increased the efficiency of triphosphate incorporation, we explored the analogs d5FM and dNaM (Figure 9B).<sup>31</sup> Gratifyingly, opposite d5SICS, both d5FMTP and dNaMTP are more efficiently inserted than dMMO2TP. Although the former is limited by the synthesis and extension of the dA–d5FM mispair, the efficiency and fidelity of replicating DNA containing the dNaM–d5SICS UBP is excellent ( $k_{\text{cat}}/K_{\text{M}} = 5.0 \times 10^6 \text{ M}^{-1} \text{ min}^{-1}$  and  $3.7 \times 10^7 \text{ M}^{-1} \text{ min}^{-1}$ , for insertion of dNaMTP opposite d5SICS and d5SICSTP opposite dNaM, respectively, and  $1.2 \times 10^6 \text{ M}^{-1} \text{ min}^{-1}$  and  $2.7 \times 10^6 \text{ M}^{-1} \text{ min}^{-1}$ , respectively, for each subsequent extension). In fact, dNaM–d5SICS was the first of our UBP candidates to be amplified in a standard PCR reaction with high efficiency and fidelity.<sup>32</sup> To determine whether dNaM is optimal for pairing with d5SICS, we further explored derivatization of the dMMO2 scaffold with Cl, Br, I, Me, and Pr meta-substituents, but none improved replication.<sup>33</sup> We next examined 28 analogs with different para-substituents, and the SAR was augmented



**Figure 9.** (A) Para- and meta-substituted dMMO2 analogs. (B) Optimized dMMO2 analogs. Sugar and phosphate groups omitted for clarity.

with PCR amplification assays using Taq, OneTaq, and KOD polymerases (Figure 9).<sup>34,35</sup> Several of the most promising analogs were further derivatized with a meta-fluoro or methoxy substituent. These efforts identified dEMO, dFIMO, and dFEMO as better partners for d5SICS than dMMO2 (Figure 9B); however, none was more optimal in these *in vitro* assays than dNaM.

## PARTS STANDARDIZATION

The successful development of a synthetic biology “part” includes not only its optimization for function, but also its standardization for function in different contexts, which here corresponds to recognition of the unnatural triphosphates by different DNA polymerases and within different local sequence contexts. To explore the extent of standardization of dNaMTP and d5SICSTP, we first examined the fidelity with which DNA containing the corresponding UBP was PCR-amplified in a

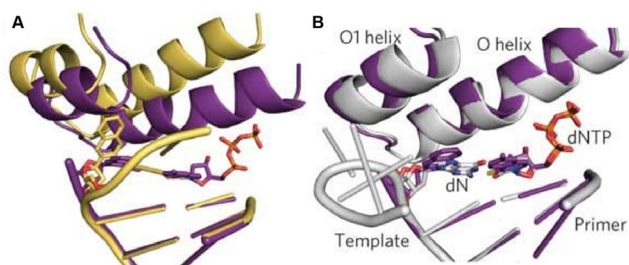
variety of different sequence contexts using DeepVent, Taq, or Phusion polymerase, which demonstrated that retention of the UBP was near or in excess of 99% per doubling.<sup>36</sup> DNA containing a UBP within a 40-nt stretch of randomized natural nucleotides was PCR amplified  $10^{24}$ -fold with OneTaq, and the products were analyzed by deep sequencing. A slight enrichment of the sequence 5'-GCNaM was observed, but this bias is no greater than that observed with natural sequences, demonstrating a sufficient level of standardization of the dNaM–d5SICS UBP for *in vivo* use.

## A SECOND STRUCTURAL INTERLUDE

Having identified a family of promising UBPs, we again returned to a structural characterization. Collaborating with the Dwyer group, we first solved the structure of a DNA duplex containing dMMO2–d5SICS. Surprisingly, the UBP still formed a cross-strand intercalated structure,<sup>37</sup> and solution-

state NOEs revealed that the dNaM–dSSICS UBP did as well (Figure 5B),<sup>38</sup> although significantly less so than the dPICS self-pair. While our interests centered around function and not structure, this mode of pairing with dNaM–dSSICS raised a perplexing question: because the UBP resembles a natural mispair more than a correct pair, how is it efficiently recognized by DNA polymerases, which are known to have evolved to reject triphosphates that form mispairs?

To address this question, we collaborated with the Marx group to solve the structures of the binary complex of KlenTaq DNA polymerase bound to a primer–template with dNaM in the templating position and with a primer terminating with a ddC, as well as the corresponding ternary complex with dSSICSTP bound. The data revealed that formation of the UBP drives the same large conformational change of the polymerase caused by the formation of a natural base pair (Figure 10)<sup>39,40</sup> and, remarkably, that the conformational

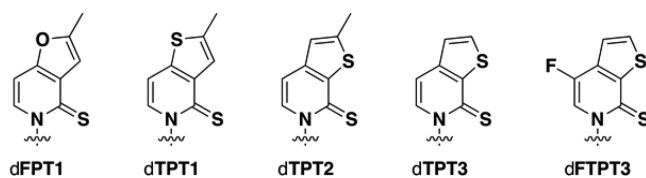


**Figure 10.** (A) Superimposition of the binary complex of KlenTaq polymerase bound to DNA with dNaM in the templating position in the open conformation (yellow) and the corresponding ternary complex bound to dSSICSTP in the closed conformation (purple). (B) Superimposition of ternary complex between KlenTaq polymerase, dNaM template DNA, and dSSICSTP (purple), or a natural dG template and dCTP (gray). Reproduced from ref 38. Copyright 2012 Nature Publishing Group.

change of the polymerase reciprocally drives a conformational change within the UBP, causing it to adopt an edge-to-edge paired natural Watson–Crick-like structure (Figure 5C). Thus, while a natural base pair is replicated with an induced-fit mechanism, the UBP is replicated with a similar, but subtly different, mutually induced-fit mechanism, providing a mechanistic explanation for its efficient replication. Nonetheless, after synthesis, the nascent UBP again adopts a cross-strand intercalated structure,<sup>41</sup> explaining the SAR data that identified a requirement for deintercalation prior to extension. At this point, we surmised that any further optimization of the UBP would require the careful optimization of intrastrand packing with neighboring natural nucleotides over cross-strand intercalation.

### ■ GROWING THE FAMILY OF *IN VITRO* REPLICATED UBPs

Based on the proposed mechanism of replication, we speculated that dNaM–dSSICS might be optimized by distal ring contraction and heteroatom derivatization of dSSICS, potentially favoring deintercalation, and thus extension, while preserving synthesis efficiency. Thus, we explored four derivatives with the distal ring replaced with para- or meta-linked thienyl, methyl furanyl, or methyl thienyl rings, as well as an additional derivative to explore fluorination at the meta position (Figure 11).<sup>42</sup> Gratifyingly, we found that dTPT2,



**Figure 11.** Distal ring-contracted dSSICS analogs. Sugar and phosphate groups omitted for clarity.

dTPT3, and dFTPT3 form UBPs with dNaM that are more efficiently replicated within DNA than the dNaM–dSSICS UBP (as demonstrated by a pre-steady-state kinetic assay, as the steady-state rates were now limited by product dissociation<sup>43</sup>). The most efficiently replicated was dNaM–dTPT3, which thus emerged as our lead UBP.

At this point, we had explored several analogs of both dMMO2 and dSSICS since the identification of dMMO2–dSSICS, but we had never explored these analogs as partners with previous generation analogs. Thus, using a PCR-based screen, we examined the amplification of DNA containing 111 different unnatural nucleotides (resulting in approximately 6000 candidate UBPs) drawn from our now expanded set of analogs. While we found that dNaM–dTPT3 is generally the most efficiently replicated of the UBPs examined, we identified seven additional and physicochemically distinct UBPs that are replicated significantly better than dNaM–dSSICS (Figure 12). Again drawing on established tenets of medicinal chemistry, these results are important, because our long-term goal of using the UBPs in a living SSO brings with it additional constraints, which may be differently satisfied by UBPs with differing physicochemical properties.

### ■ *IN VIVO* PERFORMANCE AND OPTIMIZATION

In 2014, we demonstrated that when dNaMTP and dSSICSTP are imported into *Escherichia coli* via transgenic expression of an algal nucleoside triphosphate transporter, they are used by cellular polymerases to replicate DNA containing the UBP.<sup>5</sup> However, unlike *in vitro* replication, replication in this SSO showed significant sequence biases, some of which were still observed with the dNaM–dTPT3 UBP.<sup>6</sup> This is perhaps not surprising considering that the *in vivo* replication environment is distinct from the *in vitro* environment, which had been used to generate the SAR that drove UBP optimization and standardization. Thus, we screened 135 candidate UBPs, drawn from 91 unnatural triphosphates selected to cover the range of analogs that had been explored *in vitro*, for those that when added to the media were able to support high level retention of the corresponding UBP on a plasmid within the SSO.<sup>44</sup> Much of the SAR generated was consistent with that generated *in vitro*, but there were several interesting differences. In particular, the *in vivo* environment was somewhat more permissive to the nature of the ortho substituent that had proven so critical for *in vitro* replication (although the best UBPs still retained these H-bond acceptors). Perhaps more importantly, the *in vivo* screen identified four additional UBPs, each formed by pairing a dNaM analog opposite dTPT3, that are more efficiently replicated and with less sequence bias than dNaM–dTPT3 (Figure 13), the most promising of which was dCNMO–dTPT3. While this UBP is at present the most promising lead for further SSO development, it is again the physicochemical diversity offered by a family of well replicated UBPs that is likely to prove most valuable as our efforts shift

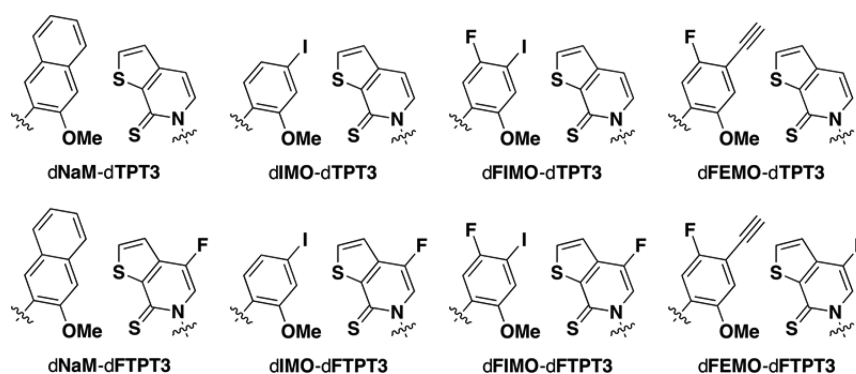


Figure 12. A family of well replicated dNaM–dTPT3-like UBPs. Sugar and phosphate groups omitted for clarity.

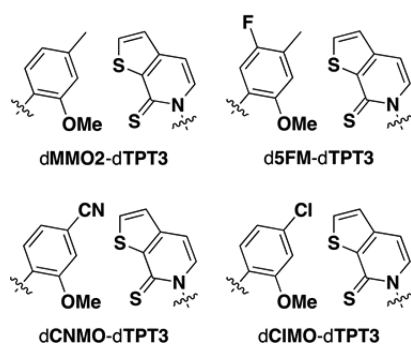


Figure 13. A family of UBPs optimized for *in vivo* expansion of the genetic alphabet. Sugar and phosphate groups omitted for clarity.

toward achieving *in vivo* transcription and translation, which we have only just begun to explore with dNaM–dTPT3.

## CONCLUSIONS AND THE CHEMIST'S APPROACH TO SYNTHETIC BIOLOGY

We have used synthetic chemistry, coupled with the methods of medicinal chemistry, to develop a family of UBPs that function not only *in vitro* but also *in vivo* and have used them to create SSOs that can store more information in their DNA. The SARs elucidated from the examination of over 150 unnatural nucleotides have guided development and identified key elements that the unnatural nucleotides must possess, most clearly, an ortho substituent that is capable of providing a hydrophobic surface as well as an H-bond acceptor and a nucleobase surface that favors intrastrand packing over cross-strand intercalation. While this Account has recounted our efforts to optimize replication, we have also recently demonstrated that DNA containing the UBPs may be transcribed into RNA in an SSO and used during translation at the ribosome to produce proteins with noncanonical amino acids.<sup>7</sup> This lays the foundation for the creation of SSOs with forms and functions not available to their natural counterparts and thereby achieves a long-standing goal of synthetic biology.

At its core, synthetic biology aims to create parts that function within living cells, imparting them with novel attributes. While the tenets of the field were originally implicitly founded on the use of chemistry to create those parts, its modern incarnation has focused on the use of parts assembled from natural components or components intended to mimic their natural counterparts. This would seem justified by the eons of evolution that optimized the natural components for functioning in a cell, at least for a similar function. However, most natural components are recognized by or interact with

multiple other components in every cell, possibly in unknown ways, and thus their introduction may have unintended consequences. While truly synthetic parts made by chemists do not benefit from eons of evolution, they are foreign to cells, possibly even drawing upon forces not used by their natural counterparts, and thus they may be more orthogonal and possibly easier to introduce and optimize without perturbation. In this case, optimization must proceed with less information, because less is known about how the parts might interact with their biological targets, and the methodology and lessons of medicinal chemistry, which conceptually face essentially the same challenges, provide the blueprint for success. If the combination of synthetic and medicinal-like chemistry can produce molecules that function alongside those evolved by nature for the most central of its processes, to store and retrieve information, then this approach is likely to be capable of discovering molecules that effectively participate in any biological process, potentially opening a new vista for chemists in synthetic biology.

## AUTHOR INFORMATION

### Corresponding Author

\*E-mail: [floyd@scripps.edu](mailto:floyd@scripps.edu).

### ORCID

Aaron W. Feldman: 0000-0002-9495-2357

Floyd E. Romesberg: 0000-0001-6317-1315

### Author Contributions

The manuscript was written through contributions of all authors. All authors have given approval to the final version of the manuscript.

### Funding

Financial support was provided by the National Institutes of Health (GM060005 and GM118178 to F.E.R.) and the National Science Foundation Graduate Research Fellowship (DGE-1346837 to A.W.F.)

### Notes

The authors declare the following competing financial interest(s): F.E.R. has a financial interest (shares) in Synthorx Inc., a company that has commercial interests in the UBP. The other author declares no other competing financial interests.

### Biographies

**Aaron W. Feldman** received his A.B. in Chemistry from Connecticut College in 2013. He is currently a graduate student in the Romesberg



lab working to chemically optimize semi-synthetic organisms with expanded genetic alphabets.

**Floyd E. Romesberg** received his Ph.D. in Chemistry from Cornell University in 1995. He is currently Professor of Chemistry at The Scripps Research Institute.

## ACKNOWLEDGMENTS

We thank former and current members of the Romesberg group, as well as those of our collaborators in the Dwyer and Marx groups for their contributions to the research highlighted in this Account.

## REFERENCES

- (1) Leduc, S. *The Mechanisms of Life*; Rebman Company: New York, 1911.
- (2) Benner, S. A.; Karalkar, N. B.; Hoshika, S.; Laos, R.; Shaw, R. W.; Matsuura, M.; Fajardo, D.; Moussatche, P. Alternative Watson-Crick synthetic genetic systems. *Cold Spring Harbor Perspect. Biol.* **2016**, *8*, a023770.
- (3) Moran, S.; Ren, R. X.; Kool, E. T. A thymidine triphosphate shape analog lacking Watson-Crick pairing ability is replicated with high sequence selectivity. *Proc. Natl. Acad. Sci. U. S. A.* **1997**, *94*, 10506–10511.
- (4) Hiraio, I.; Kimoto, M.; Yamashige, R. Natural versus artificial creation of base pairs in DNA: origin of nucleobases from the perspectives of unnatural base pair studies. *Acc. Chem. Res.* **2012**, *45*, 2055–2065.
- (5) Malyshev, D. A.; Dhami, K.; Lavergne, T.; Chen, T.; Dai, N.; Foster, J. M.; Correa, I. R., Jr.; Romesberg, F. E. A semi-synthetic organism with an expanded genetic alphabet. *Nature* **2014**, *509*, 385–388.
- (6) Zhang, Y.; Lamb, B. M.; Feldman, A. W.; Zhou, A. X.; Lavergne, T.; Li, L.; Romesberg, F. E. A semisynthetic organism engineered for the stable expansion of the genetic alphabet. *Proc. Natl. Acad. Sci. U. S. A.* **2017**, *114*, 1317–1322.
- (7) Zhang, Y.; Ptacin, J. L.; Fischer, E. C.; Aerni, H. R.; Caffaro, C. E.; San Jose, K.; Feldman, A. W.; Turner, C. R.; Romesberg, F. E. A semi-synthetic organism that stores and retrieves increased genetic information. *Nature* **2017**, DOI: 10.1038/nature24659.
- (8) Hwang, G. T.; Hari, Y.; Romesberg, F. E. The effects of unnatural base pairs and mismatches on DNA duplex stability and solvation. *Nucleic Acids Res.* **2009**, *37*, 4757–4763.
- (9) Petruska, J.; Goodman, M. F.; Boosalis, M. S.; Sowers, L. C.; Cheong, C.; Tinoco, I., Jr. Comparison between DNA melting thermodynamics and DNA polymerase fidelity. *Proc. Natl. Acad. Sci. U. S. A.* **1988**, *85*, 6252–6256.
- (10) Ogawa, A. K.; Wu, Y. Q.; McMinn, D. L.; Liu, J. Q.; Schultz, P. G.; Romesberg, F. E. Efforts toward the expansion of the genetic alphabet: Information storage and replication with unnatural hydrophobic base pairs. *J. Am. Chem. Soc.* **2000**, *122*, 3274–3287.
- (11) Ogawa, A. K.; Wu, Y. Q.; Berger, M.; Schultz, P. G.; Romesberg, F. E. Rational design of an unnatural base pair with increased kinetic selectivity. *J. Am. Chem. Soc.* **2000**, *122*, 8803–8804.
- (12) Berger, M.; Wu, Y.; Ogawa, A. K.; McMinn, D. L.; Schultz, P. G.; Romesberg, F. E. Universal bases for hybridization, replication and chain termination. *Nucleic Acids Res.* **2000**, *28*, 2911–2914.
- (13) McMinn, D. L.; Ogawa, A. K.; Wu, Y.; Liu, J.; Schultz, P. G.; Romesberg, F. E. Efforts toward expansion of the genetic alphabet: DNA polymerase recognition of a highly stable, self-pairing hydrophobic base. *J. Am. Chem. Soc.* **1999**, *121*, 11585–11586.
- (14) Yu, C.; Henry, A. A.; Romesberg, F. E.; Schultz, P. G. Polymerase recognition of unnatural base pairs. *Angew. Chem., Int. Ed.* **2002**, *41*, 3841–3844.
- (15) Henry, A. A.; Yu, C.; Romesberg, F. E. Determinants of unnatural nucleobase stability and polymerase recognition. *J. Am. Chem. Soc.* **2003**, *125*, 9638–9646.
- (16) Wu, Y.; Ogawa, A. K.; Berger, M.; McMinn, D. L.; Schultz, P. G.; Romesberg, F. E. Efforts toward expansion of the genetic alphabet: optimization of interbase hydrophobic interactions. *J. Am. Chem. Soc.* **2000**, *122*, 7621–7632.
- (17) Matsuda, S.; Fillo, J. D.; Henry, A. A.; Rai, P.; Wilkens, S. J.; Dwyer, T. J.; Geierstanger, B. H.; Wemmer, D. E.; Schultz, P. G.; Spraggon, G.; Romesberg, F. E. Efforts toward expansion of the genetic alphabet: Structure and replication of unnatural base pairs. *J. Am. Chem. Soc.* **2007**, *129*, 10466–10473.
- (18) Malyshev, D. A.; Romesberg, F. E. The expanded genetic alphabet. *Angew. Chem., Int. Ed.* **2015**, *54*, 11930–11944.
- (19) Matsuda, S.; Henry, A. A.; Romesberg, F. E. Optimization of unnatural base pair packing for polymerase recognition. *J. Am. Chem. Soc.* **2006**, *128*, 6369–6375.
- (20) Hwang, G. T.; Romesberg, F. E. Substituent effects on the pairing and polymerase recognition of simple unnatural base pairs. *Nucleic Acids Res.* **2006**, *34*, 2037–2045.
- (21) Henry, A. A.; Olsen, A. G.; Matsuda, S.; Yu, C. Z.; Geierstanger, B. H.; Romesberg, F. E. Efforts to expand the genetic alphabet: Identification of a replicable unnatural DNA self-pair. *J. Am. Chem. Soc.* **2004**, *126*, 6923–6931.
- (22) Matsuda, S.; Leconte, A. M.; Romesberg, F. E. Minor groove hydrogen bonds and the replication of unnatural base pairs. *J. Am. Chem. Soc.* **2007**, *129*, 5551–5557.
- (23) Spratt, T. E. Identification of hydrogen bonds between *Escherichia coli* DNA polymerase I (Klenow fragment) and the minor groove of DNA by amino acid substitution of the polymerase and atomic substitution of the DNA. *Biochemistry* **2001**, *40*, 2647–2652.
- (24) Morales, J. C.; Kool, E. T. Minor groove interactions between polymerase and DNA: more essential to replication than Watson-Crick hydrogen bonds? *J. Am. Chem. Soc.* **1999**, *121*, 2323–2324.
- (25) Meyer, A. S.; Blandino, M.; Spratt, T. E. *Escherichia coli* DNA polymerase I (Klenow fragment) uses a hydrogen-bonding fork from Arg668 to the primer terminus and incoming deoxynucleotide triphosphate to catalyze DNA replication. *J. Biol. Chem.* **2004**, *279*, 33043–33046.
- (26) Leconte, A. M.; Matsuda, S.; Hwang, G. T.; Romesberg, F. E. Efforts towards expansion of the genetic alphabet: pyridone and methyl pyridone nucleobases. *Angew. Chem., Int. Ed.* **2006**, *45*, 4326–4329.
- (27) Hwang, G. T.; Leconte, A. M.; Romesberg, F. E. Polymerase recognition and stability of fluoro-substituted pyridone nucleobase analogues. *ChemBioChem* **2007**, *8*, 1606–1611.
- (28) Kim, Y.; Leconte, A. M.; Hari, Y.; Romesberg, F. E. Stability and polymerase recognition of pyridine nucleobase analogues: Role of minor-groove H-bond acceptors. *Angew. Chem., Int. Ed.* **2006**, *45*, 7809–7812.
- (29) Leconte, A. M.; Hwang, G. T.; Matsuda, S.; Capek, P.; Hari, Y.; Romesberg, F. E. Discovery, characterization, and optimization of an unnatural base pair for expansion of the genetic alphabet. *J. Am. Chem. Soc.* **2008**, *130*, 2336–2343.
- (30) Hari, Y.; Hwang, G. T.; Leconte, A. M.; Joubert, N.; Heczek, M.; Romesberg, F. E. Optimization of the pyridyl nucleobase scaffold for polymerase recognition and unnatural base pair replication. *ChemBioChem* **2008**, *9*, 2796–2799.
- (31) Seo, Y. J.; Hwang, G. T.; Ordoukhanian, P.; Romesberg, F. E. Optimization of an unnatural base pair toward natural-like replication. *J. Am. Chem. Soc.* **2009**, *131*, 3246–3252.
- (32) Malyshev, D. A.; Dhami, K.; Quach, H. T.; Lavergne, T.; Ordoukhanian, P.; Torkamani, A.; Romesberg, F. E. Efficient and sequence-independent replication of DNA containing a third base pair establishes a functional six-letter genetic alphabet. *Proc. Natl. Acad. Sci. U. S. A.* **2012**, *109*, 12005–12010.
- (33) Seo, Y. J.; Romesberg, F. E. Major groove derivatization of an unnatural base pair. *ChemBioChem* **2009**, *10*, 2394–2400.
- (34) Lavergne, T.; Malyshev, D. A.; Romesberg, F. E. Major groove substituents and polymerase recognition of a class of predominantly

hydrophobic unnatural base pairs. *Chem. - Eur. J.* **2012**, *18*, 1231–1239.

(35) Lavergne, T.; Degardin, M.; Malyshev, D. A.; Quach, H. T.; Dhami, K.; Ordoukhanian, P.; Romesberg, F. E. Expanding the scope of replicable unnatural DNA: stepwise optimization of a predominantly hydrophobic base pair. *J. Am. Chem. Soc.* **2013**, *135*, 5408–5419.

(36) Malyshev, D. A.; Seo, Y. J.; Ordoukhanian, P.; Romesberg, F. E. PCR with an expanded genetic alphabet. *J. Am. Chem. Soc.* **2009**, *131*, 14620–14621.

(37) Malyshev, D. A.; Pfaff, D. A.; Ippoliti, S. I.; Hwang, G. T.; Dwyer, T. J.; Romesberg, F. E. Solution structure, mechanism of replication, and optimization of an unnatural base pair. *Chem. - Eur. J.* **2010**, *16*, 12650–12659.

(38) Betz, K.; Malyshev, D. A.; Lavergne, T.; Welte, W.; Diederichs, K.; Dwyer, T. J.; Ordoukhanian, P.; Romesberg, F. E.; Marx, A. KlenTaq polymerase replicates unnatural base pairs by inducing a Watson-Crick geometry. *Nat. Chem. Biol.* **2012**, *8*, 612–614.

(39) Rothwell, P. J.; Waksman, G. Structure and mechanism of DNA polymerases. *Adv. Protein Chem.* **2005**, *71*, 401–440.

(40) Wu, E. Y.; Beese, L. S. The structure of a high fidelity DNA polymerase bound to a mismatched nucleotide reveals an "ajar" intermediate conformation in the nucleotide selection mechanism. *J. Biol. Chem.* **2011**, *286*, 19758–19767.

(41) Betz, K.; Malyshev, D. A.; Lavergne, T.; Welte, W.; Diederichs, K.; Romesberg, F. E.; Marx, A. Structural insights into DNA replication without hydrogen bonds. *J. Am. Chem. Soc.* **2013**, *135*, 18637–18643.

(42) Li, L.; Degardin, M.; Lavergne, T.; Malyshev, D. A.; Dhami, K.; Ordoukhanian, P.; Romesberg, F. E. Natural-like replication of an unnatural base pair for the expansion of the genetic alphabet and biotechnology applications. *J. Am. Chem. Soc.* **2014**, *136*, 826–829.

(43) Morris, S. E.; Feldman, A. W.; Romesberg, F. E. Synthetic biology parts for the storage of increased genetic information in cells. *ACS Synth. Biol.* **2017**, *6*, 1834–1840.

(44) Feldman, A. W.; Dien, V. T.; Romesberg, F. E. Chemical stabilization of unnatural nucleotide triphosphates for the in vivo expansion of the genetic alphabet. *J. Am. Chem. Soc.* **2017**, *139*, 2464–2467.