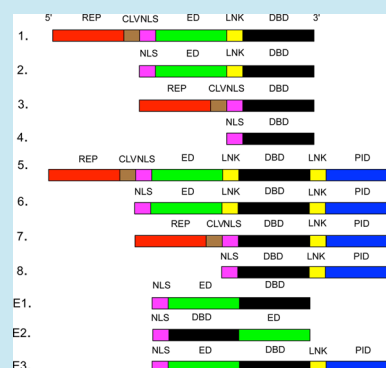# Rule-Based Design of Synthetic Transcription Factors in Eukaryotes

Oliver Purcell,[†,‡] Jean Peccoud,[§] and Timothy K. Lu*,[†,‡]

[†]Department of Electrical Engineering & Computer Science and Department of Biological Engineering, Massachusetts Institute of Technology, 77 Massachusetts Avenue, Cambridge, Massachusetts 02139, United States of America

[‡]MIT Synthetic Biology Center, 500 Technology Square, Cambridge, Massachusetts 02139, United States of America

[§]Virginia Bioinformatics Institute, Virginia Tech, Blacksburg, Virginia 24061, United States of America

**ABSTRACT:** To design and build living systems, synthetic biologists have at their disposal an increasingly large library of naturally derived and synthetic parts. These parts must be combined together in particular orders, orientations, and spacings to achieve desired functionalities. These structural constraints can be viewed as grammatical rules describing how to assemble parts together into larger functional units. Here, we develop a grammar for the design of synthetic transcription factors (sTFs) in eukaryotic cells and implement it within GenoCAD, a Computer-Aided Design (CAD) software for synthetic biology. Knowledge derived from experimental evidence was captured in this grammar to guide the user to create designer transcription factors that should operate as intended. The grammar can be easily updated and refined as our experience with using sTFs in different contexts increases. In combination with grammars that define other synthetic systems, we anticipate that this work will enable the more reliable, efficient, and automated design of synthetic cells with rich functionalities.



**KEYWORDS:** GenoCAD, grammar, synthetic transcription factor, eukaryotic transcription factor, Saccharomyces cerevisiae

Synthetic biology aims to rationally create living systems for basic science, biomedical, and biotechnology applications. To do this, one must first understand how to design synthetic networks with which to program these living systems. In order to implement complex synthetic networks, synthetic biologists require a library of well-characterized parts, such as promoters, terminators, transcription factors (TFs), and reporters, as well as rules for assembling these parts into higher-order circuits.

Transcription factors are an important class of parts for synthetic biology. They often form the regulatory links within the networks that synthetic biologists build. Synthetic networks constructed to date have largely relied upon the use of TFs from nature, such as TetR, LacI, and AraC.[1−3] However, the number of well-characterized and orthogonal natural TFs is limited; solely relying on natural TFs therefore imposes limitations on the size of synthetic networks that can be constructed. To overcome this problem, synthetic transcription factors (sTFs) have been created,[4−18] and a variety of platforms for implementing large libraries of sTFs have been described,[7,8,12] which remove the constraints imposed by natural TFs.

**Substructure of sTFs.** The design of synthetic transcription factors relies on the fact that proteins can be modularized and assembled into interchangeable, and generally quasi-independent protein domains. The term 'transcription factor' traditionally refers to any protein that regulates transcription by any means. However, within the context of this discussion, all TFs, whether synthetic or natural, influence gene expression by DNA binding at or near promoters and therefore require DNA-binding domains (DBDs). For instance, the zinc-finger-b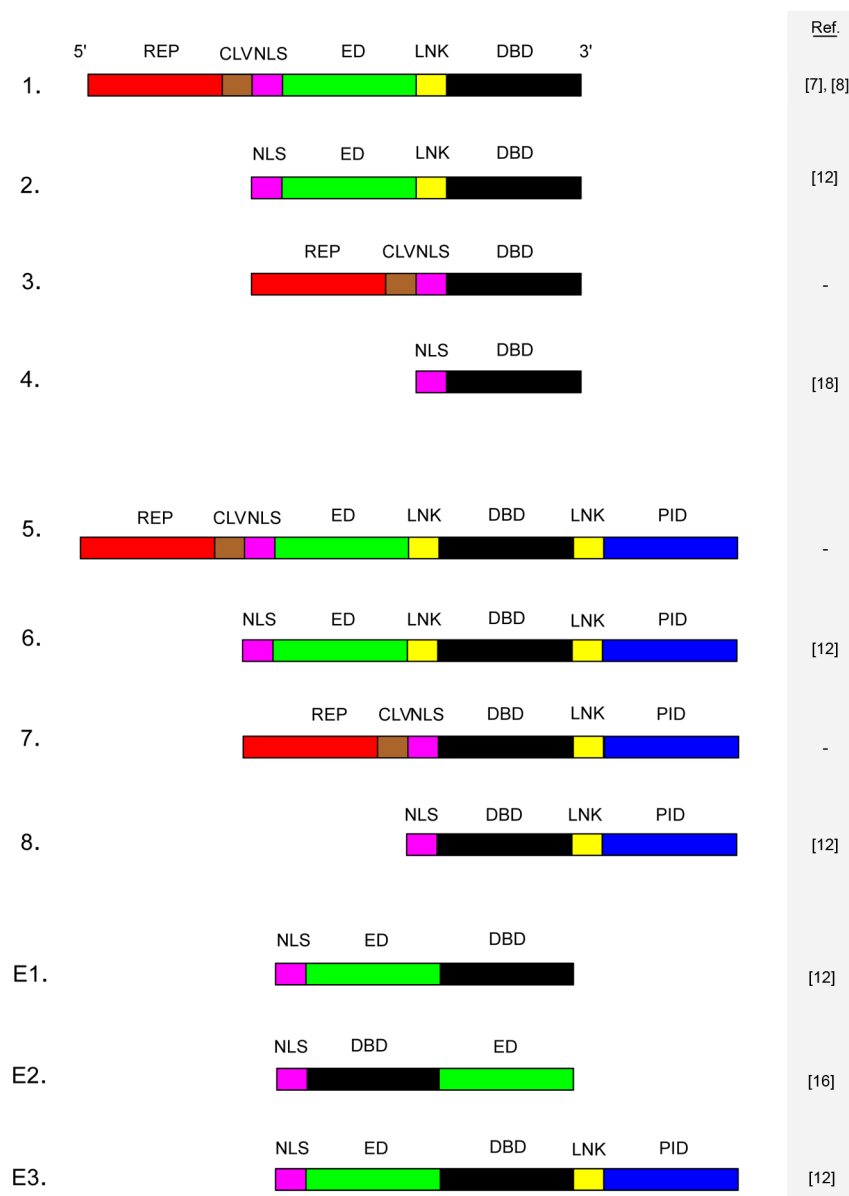ased class of sTFs uses a series of zinc fingers as DBDs, with each zinc finger (ZF) containing a defined amino-acid sequence, which recognizes a specific DNA triplet code (e.g., CTG). By fusing together multiple zinc fingers, a larger DNA-binding domain that recognizes a longer DNA sequence can be constructed.[17,19−22]

TFs can generally be divided into two classes; activators, which activate or increase transcription, and repressors, which decrease or repress transcription. In yeast, activation and repression are typically mediated by 'effector domains', which are fused to DBDs, allowing them to be targeted to specific promoters. A commonly used activation domain in synthetic activators in yeast is the VP16 domain, or its derivative, the VP64 domain (formed from 4 tandem repeats of the VP16 domain),[7,8,12] while a commonly used repression domain is the SSN6 domain.[23] VP16 recruits various transcription factors necessary for transcription and a Histone Acetylase Complex (HAC).[24] HACs lead to acetylation of nearby histones, causing chromatin to unwind and allowing access to the promoter by the transcriptional machinery.[24] Conversely, SSN6 is thought to repress transcription by preventing transcriptional initiation by RNA polymerase and recruiting Histone Deacetylase Complexes (HDACs), which deacetylate histones, leading to compaction of the chromatin and prevention of further access to the promoter by the transcriptional machinery.[25] Repression can also be achieved without an effector domain by using DBDs to sterically block initiation of RNA polymerase.[18]

**Figure 1.** Structures of sTFs allowed within the grammar. Eight possible general structures are allowed within the grammar. In addition, E1−E3 are experimentally verified structures. DBD = DNA-binding domain, LNK = linker domain, ED = effector domain, NLS = nuclear localization signal, CLV = cleavage domain, REP = reporter, PID = protein interaction domain. All constructs are oriented from 5′ to 3′. References for structures 1−8 describe studies in which similar structures have been experimentally verified. References for structures E1−E3 denote the study in which the structure was experimentally verified.

**Rules and Grammars.** A grammar is simply a set of design rules, which can be used to guide the design process or enforce standards. This formalism is suited to capturing a domain expertise in a format that constrains nonexpert users to produce designs that conform to what is known by expert users (i.e., the experienced synthetic biologist) to typically work. For instance, many synthetic biologists working on various applications use eukaryotic sTFs in their projects. Yet, only a small fraction of the potential users of sTFs are familiar enough with sTF design to take advantage of the rapid progress in this field. The grammar presented here could help transfer the expertise of sTF specialists to those with expertise in other fields.

To someone specializing in the development of the next generation of sTFs, the benefits of constraining the design process may not be immediately apparent since optimal designs are unknown. In this case, grammars are a formal representation of a hypothesis that will be tested experimentally. This formalization effort encourages a thorough analysis of the different aspects of the design process, which can help uncover potential issues before starting the experimental validation. It also supports the articulation of various context-dependencies that may affect the success of a design strategy in different conditions. Furthermore, grammars implemented within computer-aided design tools may help to organize experimental libraries and plans.

**Rules for sTF Design.** On some level, all biological parts (whether natural or synthetic) conform to certain design rules to varying degrees. For example, *E. coli* promoters usually require −10 and −35 boxes for RNA polymerase binding to initiate transcription, while proteins require a start codon from where translation is started. The structure of sTFs can also be designed to conform to certain rules. For instance, to design an

sTF that behaves as an activator, it should have a DBD fused somehow to an AD. However, just as the structure of an sTF can be more complex than a two-domain fusion, the grammar can also be more complex.

Here, we propose a grammar for the design of sTFs in *Saccharomyces cerevisiae*. We implement this in GenoCAD, a web-based synthetic-biology CAD software.[26] GenoCAD was derived from the observation that constructs used in synthetic biology could be generated by context-free grammars.[27] It is therefore a logical choice for implementing an sTF grammar. GenoCAD includes a system to create and manage libraries of user-defined parts. The GenoCAD design module provides a wizard-like interface which guides users to generate structurally valid constructs, and allows the online design workspace to be customized.[26] We propose grammars for the design of sTFs based on zinc fingers, transcription activator-like effectors (TALEs), and the recently developed Clustered Regularly Interspaced Short Palindromic Repeats (CRISPR)/Cas-based system. Our grammar covers the design of sTFs that (1) use any one of these systems, (2) use effector domains to activate or repress transcription, (3) use fluorescent reporter domains to enable quantification of sTF abundance, and (4) permit the design of sTFs that form dimeric complexes with other sTFs, which can be used to engineer cooperativity between sTF monomers.

We believe our grammar serves as a first attempt to standardize sTF design and create a foundation that can be built upon and refined as experience with designing and using sTFs grows.

## ■ sTF GRAMMAR

While it would be possible to construct an arbitrarily broad grammar that would allow an expert user to define any combination of protein domains in any order, this defeats the purpose of the grammar in productively constraining nonexpert users. Therefore, we have opted for a highly constrained grammar based around 11 possible sTF structures (Figure 1). The 11 possible structures are, 5′ to 3′, as follows:

1. 5′-REP-CLV-NLS-ED-LNK-DBD-3′
2. 5′-NLS-ED-LNK-DBD-3′
3. 5′-REP-CLV-NLS-DBD-3′
4. 5′-NLS-DBD-3′
5. 5′-REP-CLV-NLS-ED-LNK-DBD-LNK-PID-3′
6. 5′-NLS-ED-LNK-DBD-LNK-PID-3′
7. 5′-REP-CLV-NLS-DBD-LNK-PID-3′
8. 5′-NLS-DBD-LNK-PID-3′
E1. 5′-NLS-ED-DBD-3′
E2. 5′-NLS-DBD-ED-3′
E3. 5′-NLS-ED-DBD-LNK-PID-3′

where PID = protein interaction domain, LNK = linker domain, ED = effector domain, CLV = cleavage domain, REP = reporter domain, NLS = nuclear localization sequence.

Structures 1–8 shown in Figure 1 allow for the construction of sTFs that can provide either activation through effector domains, or repression by effector domains or steric hindrance of RNA polymerase initiation. The sTF expression levels can be quantified using reporter proteins and sTFs can be made to behave cooperatively when paired with a suitable partner. These structures therefore cover the range of functions that are required by sTFs in the construction of synthetic gene networks. The design of our structures is based on a synthesis of the available experimental evidence. However, many of these
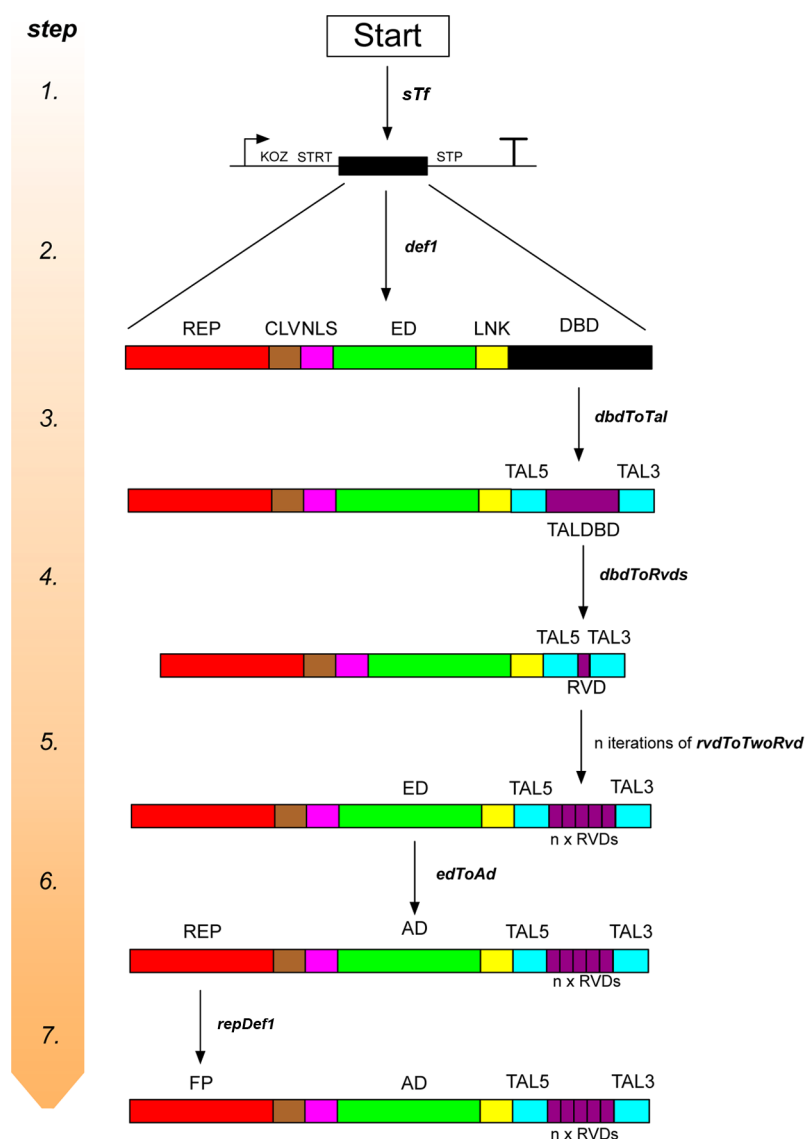
structures are themselves novel and, to our knowledge, have not yet been experimentally verified. References in Figure 1 denote studies that offer experimental evidence for structures that are similar to the structures presented here. We have also included three structures that have been experimentally verified in *S. cerevisiae* (E1–E3). These are two variants of an NLS-DBD-ED structure (E1 and E2) and a modification of E1 that allows for protein–protein interactions (E3).

For structures 1–8, the general structural constraints captured in the grammar are as follows:

- The physical structure of the sTF (i.e., the ordering of the domains) is organized around the position of the DBD.
- All domains apart from the PID (with a LNK domain present between it and the DBD) are built 5′ to the DBD.
- PIDs are built 3′ to the DBD.
- If a reporter is used, it is the 5′ terminal domain.
- Between a DBD and either a PID or an ED, there must be a linker domain (LNK).
- Between any domain and the reporter domain, there must be a cleavage domain (CLV). The most common cleavage domain, and the one used in the library with our GenoCAD grammar is the 2A sequence. However, this is not a true cleavage domain as no proteolytic cleavage of the protein occurs. Rather, the 2A sequence causes 'ribosome skipping' to occur,[28] whereby the peptide bond formation does not occur, and two separate proteins are therefore produced.[29] In order to follow the literature on 2A domains, we refer to CLVs as 'cleavage domains'[29] to denote that they include any domain that separates proteins, whether it be by true cleavage or not.
- A nuclear localization signal (NLS) is added at the 5′ of the protein. If a cleavage domain is present, then the NLS is immediately 3′ to the cleavage domain.

There are many possible variations of the structures 1–8 shown in Figure 1, and the subdomains in these different fusion configurations may have different structures and therefore different activities. For simple fusion proteins, such as the fusion of one domain with a fluorescent protein, it has been suggested that both configurations of the fusion protein be tested.[30] However, when the number of possible configurations is large, testing all possible configurations is usually impractical. The justifications for these general constraints are as follows:

- Both the PDZ and leucine zipper (LZ) domains have been used successfully as PIDs to enable cooperativity in sTFs.[8,12] LZ domains have been shown to function when placed internally in the sTF[8] and should also function at either terminus. However, the ligand to which the PDZ domain binds must be at the C-terminal of the protein.[31] To minimize the number of available structures, we therefore constrain both components of a PID-based interaction (the protein and its ligand) to be at the 3′ (C-terminal) end.
- Because of the constraint on the PID to be at the 3′ end, we therefore constrain all other domains to be 5′ to the DBD. ZF-based sTFs have been successfully constructed with the effector domain to the 5′ of the DBD.[12]
- Linker domains are routinely used when creating fusion proteins and have been shown to improve folding and stability of fusion proteins, enhance the expression of

**Figure 2.** Example design process for an sTF within GenoCAD. The seven steps of the process are oriented from top to bottom. The transformation rules that transform the construct from the start state to the final construct are depicted in bold italics (sTf, def1, dbdToTal, dbdToRvds, rvdToTwoRvd, edToAd, repDef1). DBD = DNA-binding domain, LNK = linker domain, ED = effector domain, NLS = nuclear localization signal, CLV = cleavage domain, REP = reporter, TALBDB = TALE DNA-binding domain, TAL5 = 5′ domain of the TALE, TAL3 = 3′ domain of the TALE, RVD = repeat variable domain, AD = activation domain, FP = fluorescent protein. The right-angled arrow and "T" denote the promoter and the terminator, respectively. KOZ, STRT, and STP denote a Kozak sequence, start codon, and stop codon, respectively. All constructs are oriented from 5′ to 3′.

fusion proteins, and increase the activity of the fusion protein.[32]

- Nuclear localization signals have been placed internal to sTFs[7] as well as at the termini.[7,12] To our knowledge, there has not yet been a comprehensive study as to if and how the placement and number of NLSs affects the characteristics of the sTF. Therefore, here, we place the NLS at the N-terminal region with respect to the DBD.

- The addition of a reporter domain at the 5′ end of the sTF allows for the concentration of sTFs present within the cell at any given time to be quantified. However, the presence of this additional reporter domain may adversely affect the folding of the rest of the sTF, impairing its function and vice versa. Placing a 'cleavage' domain before the reporter may mitigate any such issue. Upon translation, the 'cleavage' domain results in the

protein sequence being cleaved at a specific position. The efficiency of 'cleavage' with 2A domains has been shown to be affected by the sequence of the upstream protein.[29] However, by simply adding a Gly-Ser-Gly (GSG linker) before the 2A sequence, the efficiency can be increased to ∼100% for all upstream proteins tested.[29] We include this GSG linker as a standard component of the P2A sequence.

**PIDs.** PIDs can be defined as either homodimerization (e.g., the LZ domain) or heterodimerization domains (e.g., the PDZ domain and its ligand, or heterospecific interactions based on coiled-coils[33,34]). Heterodimeric PIDs can be further defined as positive (+) or negative (−). A positive PID is intended to interact with a complementary negative PID. As both PIDs are on the 3′ end of the sTF, in the case of PDZ domains one of the sTFs must reverse its direction to allow for interaction

between the two sTFs. It is trivial to synthesize the DNA-binding site of the DBD in reverse, so this is where the orientation issue is dealt with, rather than allowing PIDs to be at the 5′ end of the sTF.

**DBDs.** ZF domains cannot be further subdefined in the grammar. This is because of apparent interdependencies between the individual ZFs that form a ZF-array DBD, which means that ZFs do not always behave in a truly modular fashion.[35−37] It is therefore more reliable to use entire ZF-array DBDs that have been verified for specificity, rather than construct them *de novo* and risk interdependency issues. No such interdependencies are known for the Repeat Variable Domains (RVDs) that form the TALE-based DBDs, and therefore, TALEs can be further subdefined in the design process. A TALE domain must include a 5′ and 3′ TALE region, and >0 repeat variable domains (RVDs) in-between the 5′ and 3′ TALE DBD region. The dCas9 domain cannot be further subdefined.

*gRNA.* The CRISPR-TF system comprises a dCas9 domain (optionally fused to an effector domain) and a guide RNA (gRNA). dCas9 is a catalytically inactive form of the Cas9 nuclease. The gRNA itself is comprised of a sequence that binds through complementary base pairing to one strand of the DNA target sequence, and a 'handle' sequence: a hairpin forming sequence that dCas9 recognizes and binds. The gRNA therefore 'guides' the dCas9 based TF to its target site and determines the DNA-binding specificity of the dCas9:gRNA complex, and therefore, its effects on the expression of the target gene.[10,11,16] Every dCas9 domain should have a gRNA defined for it. We use a single gRNA, where the handle and the targeting sequence are fused, rather than the original 'dual' RNA system, where these components were separate and had to interact *in vivo* for the system to function.[38,39]

The user is able to define gRNAs within the sTF grammar. As gRNAs are not translated, they do not require either a start or stop codon.

**EDs.** Effector domains can be designated as either activator domains (ADs) or repressor domains (RDs).

**REPs.** Reporter domains can either be designated as a fluorescent protein reporter (e.g., GFP or mCherry) or a nonfluorescent protein reporter (e.g., β-galactosidase)

## ■ GENOCAD IMPLEMENTATION

The preceding section described the biological details of the grammar. This section describes the specifics of the implementation of this grammar within GenoCAD.

A GenoCAD grammar is defined by categories of genetic parts and transformation rules between these categories. For instance, an ED would be a category, as would an AD and an RD. The transformation rule that links these would be that an ED can be defined as (or 'transformed' into) either an AD or an RD. When a user wants to define a genetic construct within GenoCAD, they always begin from the 'start' category. From the start category, they can iteratively transform particular categories into different subcategories, therefore defining the specifics of the genetic construct. An illustrative example is shown in Figure 2.

The categories and transformation rules for the sTF grammar are given in Tables 1 and 2, respectively.

**Future Developments to the Grammar.** In this paper, we have presented a grammar for the design of synthetic transcription factors. We have implemented this in GenoCAD, a CAD software that uses grammars to define synthetic

### Table 1. Categories in the sTF Grammar[a]

| category ID | description | type |
|---|---|---|
| S (start) | start category; the default 'root' category of the grammar | rewritable |
| sTF | the entire sTF (not including the promoter or terminator) | rewritable |
| DBD | DBD of the sTF | rewritable |
| TALDBD | the part of the TALE formed by RVDs; does not include the 5′ and 3′ TALE ends | rewritable |
| RVD | an individual RVD that forms part of the TALE DBD | rewritable |
| ED | effector domain; can be either an activation or repression domain | rewritable |
| REP | reporter domain | rewritable |
| PID | protein interaction domain; can interact with other protein interaction domains to allow the sTF to form dimers | rewritable |
| gRNA | guide RNA | rewritable |
| PROM | promoter that drives expression of the sTF | terminal |
| KOZ | Kozak sequence | terminal |
| TERM | terminator for the sTF | terminal |
| ZFDBD | DBD for a ZF-based sTF | terminal |
| DCAS9 | catalytically inactive Cas9 domain | terminal |
| TAL5 | 5′ end of the TALE | terminal |
| TAL3 | 3′ end of the TALE | terminal |
| LNK | (usually) short linker sequence that joins two domains | terminal |
| CLV | amino acid sequence that joins two domains but is 'cleaved' during/after translation, separating the domains | terminal |
| FP | fluorescent protein that acts as a reporter | terminal |
| xREP | any domain that acts as a reporter but is not an FP | terminal |
| AD | effector domain that is an activation domain; it causes an increase in the expression of the target promoter | terminal |
| RD | effector domain that is a repression domain; it causes a decrease in the expression of the target promoter | terminal |
| PIDhm | homodimerizing PID domain | terminal |
| PIDht+ | heterodimerizing 'positive' PID domain; will interact with (bind to) its corresponding PIDht− domain | terminal |
| PIDht− | heterodimerizing 'negative' PID domain; will interact with (bind to) its corresponding PIDht+ domain | terminal |
| STRT | start codon | terminal |
| TRGT | sequence of the gRNA complementary to the target sequence | terminal |
| HNDLE | Cas9-binding domain of the gRNA | terminal |
| STP | stop codon | terminal |

[a]Re-writable categories can be transformed into other categories, while terminal categories cannot.

constructs. The grammars allow for the construction of 11 different sTF structures based on commonly used components. The DNA-binding domain of the sTF can be defined as zinc fingers, TALEs, or the dCas9 protein (which acts in concert with a gRNA to target specific DNA sequences). Our grammar also allows for the design of cooperative transcription factors through the incorporation of protein interaction domains.

The grammar presented here represents one interpretation of our current experience with sTFs. However, we make two implicit assumptions in defining a grammar: first, our grammar is focused on the domain structure of the transcription factors, while ultimately it is the amino-acid sequence of the protein that is important, as it is this sequence which defines how the protein folds and therefore how it functions. Second, although we base the selection of our 8 general structures on experimental evidence, we are extrapolating from this evidence to form the structures we described here. Thus, we assume that what has been observed in one context (e.g., the placement of NLSs in a particular sTF) will also be applicable in other sTFs.

**Table 2. Transformation Rules in the sTF Grammar**

| rule code | rule | description |
|---|---|---|
| sTf | S to PROM-KOZ-STRT-sTF-STP-TERM | converts the start state to a gene structure containing an sTF |
| def1 | sTF to REP-CLV-NLS-ED-LNK-DBD | converts the sTF to the first structure variant in list Figure 1 |
| def2 | sTF to NLS-ED-LNK-DBD | converts the sTF to the second structure variant in list Figure 1 |
| def3 | sTF to REP-CLV-NLS-DBD | converts the sTF to the third structure variant in list Figure 1 |
| def4 | sTF to NLS-DBD | converts the sTF to the fourth structure variant in list Figure 1 |
| def5 | sTF to REP-CLV-NLS-ED-LNK-DBD-LNK-PID | converts the sTF to the fifth structure variant in list Figure 1 |
| def6 | sTF to NLS-ED-LNK-DBD-LNK-PID | converts the sTF to the sixth structure variant in list Figure 1 |
| def7 | sTF to REP-CLV-NLS-DBD-LNK-PID | converts the sTF to the seventh structure variant in list Figure 1 |
| def8 | sTF to NLS-DBD-LNK-PID | converts the sTF to the eighth structure variant in list Figure 1 |
| defE1 | sTF to NLS-ED-DBD | converts the sTF to the structure E1 in Figure 1 |
| defE2 | sTF to NLS-DBD-ED | converts the sTF to the structure E2 in Figure 1 |
| defE3 | sTF to NLS-ED-DBD-LNK-PID | converts the sTF to the structure E3 in Figure 1. |
| dbdToTal | DBD to TAL5-TALDBD-TAL3 | converts the DBD to a TALE DBD including the 5′ and 3′ TALE end regions |
| dbdToZf | DBD to ZF | converts the DBD to a zinc finger |
| dbdToDcas9 | DBD to DCAS9 | converts the DBD to a dCas9 domain |
| edToAd | ED to AD | converts the effector domain to an activation domain |
| edToRd | ED to RD | converts the effector domain to a repression domain |
| dbdToRvds | TALDBD to RVD | converts the TALE DBD to an RVD |
| rvdToTwoRvd | RVD to RVD-RVD | converts one RVD domain to two RVD domains |
| repDef1 | REP to FP | converts a reporter to a fluorescent protein |
| repDef2 | REP to xREP | converts a reporter to a reporter domain other than a fluorescent protein |
| pidhmDef | PID to PIDhm | converts a PID to a PIDhm domain |
| pidht+Def | PID to PIDht+ | converts a PID to a PIDht+ domain |
| pidht-Def | PID to PIDht- | converts a PID to a PIDht− domain |
| grna | S to PROM-gRNA-TERM | converts the start state to a gene structure containing a gRNA |
| grnaDef | gRNA to TRGT-HNDLE | converts the gRNA to a target sequence and a handle sequence |

These rules are intended to allow a user to design sTFs with structures that will be functional. It should be noted that the 8 general structures we present in Figure 1 have not yet been experimentally verified for functionality—although there are similarities to known functional structures. However, there likely exist structures that will have more desirable characteristics than the ones allowed within this grammar. For instance, perhaps using multiple nuclear-localization sequences in various specific positions may increase the rate of nuclear import for a certain sTF[40,41] or putting a longer linker in between a ZF DBD and a particular ED may increase the magnitude of the expression change caused by the ED.[32] This grammar should therefore be revised as our knowledge of sTF design increases.

This grammar could be improved in a number of ways. For example, although our grammar allows for a TALE DBD to be constructed with only a single RVD, in reality, to ensure both sufficient specificity and binding affinity, the number of RVDs would typically be on the order of 20.[7] With a single rewriting rule (RVD → RVD RVD), the grammar can introduce as many RVDs as necessary. However, the process is cumbersome and having many RVD icons in the design is not particularly elegant. A more refined version of the grammar could introduce categories representing blocks of 1, 5, 10 RVDs and the corresponding rules. Future iterations of the grammar will make it possible to quickly generate a broad range of RVDs using a smaller number of icons and rewriting steps. Furthermore, the PID domains are labeled 'positive' and 'negative', which guide the user somewhat toward permissible pairings of sTFs. However, this is not a constraint, and the user is still able to pair sTFs incorrectly. An improvement would therefore be for the user to 'pair up' designed sTFs within GenoCAD, which could be automatically examined for compatibility. Another useful constraint on pairing would be between dCas9 domains and gRNAs.

The current version of the grammar focuses on the design of individual transcription factors. A natural extension of this grammar would be to include rules allowing the design of gene networks derived from these sTFs. For instance, one could constrain sTFs to 'pair' with promoters that contain sequences to which the DBD of the sTF is able to bind. Adding a network layer to the grammar would make it possible to benefit from the GenoCAD simulation environment. As sTF libraries become better characterized with kinetic data, it would be advantageous to be able to incorporate this information into GenoCAD for the purpose of simulating the dynamics of gene networks built from these sTFs. Further integration of synthetic circuit modeling within whole-cell models in GenoCAD could enhance the utility of this approach.[42]

As the number and complexity of components engineered by synthetic biologists increases, encapsulating current knowledge by defining standards will become increasingly important. These standards will allow for more reliable construction of synthetic living systems by scientists and engineers with a more wide-ranging level of expertise. We propose that sTF grammars, such as those presented here, begin to be considered as a first step toward the standardization of a broad range of synthetic genetic parts that could be combined in synthetic gene circuit designs.

## ■ AUTHOR INFORMATION

**Corresponding Author**
*E-mail: timlu@mit.edu.

**Author Contributions**
O.P. and J.P. implemented the grammars. O.P., J.P., and T.K.L. designed and analyzed the grammars and wrote the paper.

**Notes**
The authors declare the following competing financial interest(s): J.P. has a financial interest in GenoFAB, LLC.
The sTF grammar is available upon request.

## ■ REFERENCES

(1) Elowitz, M. B., and Leibler, S. (2000) A synthetic oscillatory network of transcriptional regulators. *Nature 403*, 335–338.

(2) Gardner, T. S., Cantor, C. R., and Collins, J. J. (2000) Construction of a genetic toggle switch in *Escherichia coli. Nature 403*, 339–342.

(3) Lou, C., Liu, X., Ni, M., Huang, Y., Huang, Q., Huang, L., Jiang, L., Lu, D., Wang, M., Liu, C., Chen, D., Chen, C., Chen, X., Yang, L., Ma, H., Chen, J., and Ouyang, Q. (2010) Synthesizing a novel genetic sequential logic circuit: A push-on push-off switch. *Mol. Syst. Biol. 6*, 350.

(4) Perez-Pinera, P., Kocak, D. D., Vockley, C. M., Adler, A. F., Kabadi, A. M., Polstein, L. R., Thakore, P. I., Glass, K. A., Ousterout, D. G., Leong, K. W., Guilak, F., Crawford, G. E., Reddy, T. E., and Gersbach, C. A. (2013) RNA-guided gene activation by CRISPR-Cas9-based transcription factors. *Nat. Methods 10*, 973–976.

(5) Perez-Pinera, P., Ousterout, D. G., Brunger, J. M., Farin, A. M., Glass, K. A., Guilak, F., Crawford, G. E., Hartemink, A. J., and Gersbach, C. A. (2013) Synergistic and tunable human gene activation by combinations of synthetic transcription factors. *Nat. Methods 10*, 239–242.

(6) Bogdanove, A. J., and Voytas, D. F. (2011) TAL effectors: Customizable proteins for DNA targeting. *Science 333*, 1843–1846.

(7) Garg, A., Lohmueller, J. J., Silver, P. A., and Armel, T. Z. (2012) Engineering synthetic TAL effectors with orthogonal target sites. *Nucleic Acids Res. 40*, 7584–7595.

(8) Lohmueller, J. J., Armel, T. Z., and Silver, P. A. (2012) A tunable zinc finger-based framework for Boolean logic computation in mammalian cells. *Nucleic Acids Res. 40*, 5180–5187.

(9) Zhang, F., Cong, L., Lodato, S., Kosuri, S., Church, G. M., and Arlotta, P. (2011) Efficient construction of sequence-specific TAL effectors for modulating mammalian transcription. *Nat. Biotechnol. 29*, 149–153.

(10) Bikard, D., Jiang, W., Samai, P., Hochschild, A., Zhang, F., and Marraffini, L. A. (2013) Programmable repression and activation of bacterial gene expression using an engineered CRISPR-Cas system. *Nucleic Acids Res. 41*, 7429–7437.

(11) Qi, L. S., Larson, M. H., Gilbert, L. A., Doudna, J. A., Weissman, J. S., Arkin, A. P., and Lim, W. A. (2013) Repurposing CRISPR as an RNA-guided platform for sequence-specific control of gene expression. *Cell 152*, 1173–1183.

(12) Khalil, A. S., Lu, T. K., Bashor, C. J., Ramirez, C. L., Pyenson, N. C., Joung, J. K., and Collins, J. J. (2012) A synthetic biology framework for programming eukaryotic transcription functions. *Cell 150*, 647–658.

(13) Morbitzer, R., Römer, P., Boch, J., and Lahaye, T. (2010) Regulation of selected genome loci using de novo-engineered transcription activator-like effector (TALE)-type transcription factors. *Proc. Natl. Acad. Sci. U.S.A. 107*, 21617–21622.

(14) Folcher, M., Xie, M., Spinnler, A., and Fussenegger, M. (2013) Synthetic mammalian trigger-controlled bipartite transcription factors. *Nucleic Acids Res. 41*, e134.

(15) Geissler, R., Scholze, H., Hahn, S., Streubel, J., Bonas, U., Behrens, S.-E., and Boch, J. (2011) Transcriptional activators of human genes with programmable DNA-specificity. *PLoS One 6*, e19509.

(16) Farzadfard, F., Perli, S. D., and Lu, T. K. (2013) Tunable and multifunctional eukaryotic transcription factors based on CRISPR/Cas. *ACS Synth. Biol. 2*, 604–613.

(17) Beerli, R. R., and Barbas, C. F. (2002) Engineering polydactyl zinc-finger transcription factors. *Nat. Biotechnol. 20*, 135–141.

(18) Blount, B. A., Weenink, T., Vasylechko, S., and Ellis, T. (2012) Rational diversification of a promoter providing fine-tuned expression and orthogonal regulation for synthetic biology. *PLoS One 7*, e33279.

(19) Pabo, C. O., Peisach, E., and Grant, R. A. (2001) Design and selection of novel Cys2His2 zinc finger proteins. *Annu. Rev. Biochem. 70*, 313–340.

(20) Liu, Q., Segal, D. J., Ghiara, J. B., and Barbas, C. F. (1997) Design of polydactyl zinc-finger proteins for unique addressing within complex genomes. *Proc. Natl. Acad. Sci. U.S.A. 94*, 5525–5530.

(21) Kim, J. S., and Pabo, C. O. (1998) Getting a handhold on DNA: Design of poly-zinc finger proteins with femtomolar dissociation constants. *Proc. Natl. Acad. Sci. U.S.A. 95*, 2812–2817.

(22) Sera, T., and Uranga, C. (2002) Rational design of artificial zinc-finger proteins using a nondegenerate recognition code table. *Biochemistry 41*, 7074–7081.

(23) Keleher, C. A., Redd, M. J., Schultz, J., Carlson, M., and Johnson, A. D. (1992) Ssn6-Tup1 is a general repressor of transcription in yeast. *Cell 68*, 709–719.

(24) Hall, D. B., and Struhl, K. (2002) The VP16 activation domain interacts with multiple transcriptional components as determined by protein–protein cross-linking *in vivo. J. Biol. Chem. 277*, 46043–46050.

(25) Malavé, T. M., and Dent, S. Y. R. (2006) Transcriptional repression by Tup1–Ssn6. *Biochem. Cell Biol. 84*, 437–443.

(26) Czar, M. J., Cai, Y., and Peccoud, J. (2009) Writing DNA with GenoCAD. *Nucleic Acids Res. 37*, W40–W47.

(27) Cai, Y., Hartnett, B., Gustafsson, C., and Peccoud, J. (2007) A syntactic model to design and verify synthetic genetic constructs derived from standard biological parts. *Bioinformatics 23*, 2760–2767.

(28) Donnelly, M. L. L., Luke, G., Mehrotra, A., Li, X., Hughes, L. E., Gani, D., and Ryan, M. D. (2001) Analysis of the aphthovirus 2A/2B polyprotein 'cleavage' mechanism indicates not a proteolytic reaction, but a novel translational effect: A putative ribosomal 'skip'. *J. Gen. Virol. 82*, 1013–1025.

(29) Szymczak-Workman, A. L., Vignali, K. M., and Vignali, D. A. A. (2012) Design and construction of 2A peptide-linked multicistronic vectors. *Cold Spring Harb. Protoc. 2012*, 199–204.

(30) Snapp, E. (2005) Design and use of fluorescent fusion proteins in cell biology. *Curr. Protoc. Cell Biol.*, 21.4.1–21.4.13.

(31) Harris, B. Z., and Lim, W. A. (2001) Mechanism and role of PDZ domains in signaling complex assembly. *J. Cell Sci. 114*, 3219–3231.

(32) Chen, X., Zaro, J. L., and Shen, W.-C. (2012) Fusion protein linkers: Property, design and functionality. *Adv. Drug Delivery Rev. 65*, 1357–1369.

(33) Thompson, K. E., Bashor, C. J., Lim, W. A., and Keating, A. E. (2012) SYNZIP protein interaction toolbox: *in vitro* and *in vivo* specifications of heterospecific coiled-coil interaction domains. *ACS Synth. Biol. 1*, 118–129.

(34) Reinke, A. W., Grant, R. A., and Keating, A. E. (2010) A synthetic coiled-coil interactome provides heterospecific modules for molecular engineering. *J. Am. Chem. Soc. 132*, 6025–6031.

(35) Ramirez, C. L., Foley, J. E., Wright, D. A., Müller-Lerch, F., Rahman, S. H., Cornu, T. I., Winfrey, R. J., Sander, J. D., Fu, F., Townsend, J. A., Cathomen, T., Voytas, D. F., and Joung, J. K. (2008) Unexpected failure rates for modular assembly of engineered zinc fingers. *Nat. Methods 5*, 374–375.

(36) Lam, K. N., van Bakel, H., Cote, A. G., van der Ven, A., and Hughes, T. R. (2011) Sequence specificity is obtained from the majority of modular C2H2 zinc-finger arrays. *Nucleic Acids Res. 39*, 4680–4690.

(37) Carlson, D. F., Fahrenkrug, S. C., and Hackett, P. B. (2012) Targeting DNA With Fingers and TALENs. *Mol. Ther. Nucleic Acids 1*, e3.

(38) Jinek, M., Chylinski, K., Fonfara, I., Hauer, M., Doudna, J. A., and Charpentier, E. (2012) A programmable dual-RNA-guided DNA endonuclease in adaptive bacterial immunity. *Science 337*, 816–821.

(39) Mali, P., Yang, L., Esvelt, K. M., Aach, J., Guell, M., DiCarlo, J. E., Norville, J. E., and Church, G. M. (2013) RNA-guided human genome engineering via Cas9. *Science 339*, 823–826.

(40) Luo, M., Pang, C. W. M., Gerken, A. E., and Brock, T. G. (2004) Multiple nuclear localization sequences allow modulation of 5-lipoxygenase nuclear import. *Traffic 5*, 847–854.

(41) Gassman, N. R., Clodfelter, J. E., McCauley, A. K., Bonin, K., Salsbury, F. R., and Scarpinato, K. D. (2011) Cooperative nuclear

localization sequences lend a novel role to the N-terminal region of MSH6. *PLoS One 6*, e17907.

(42) Purcell, O., Jain, B., Karr, J. R., Covert, M. W., and Lu, T. K. (2013) Towards a whole-cell modeling approach for synthetic biology. *CHAOS 23*, 25112.

**744**

dx.doi.org/10.1021/sb400134k | *ACS Synth. Biol.* 2014, 3, 737–744