

# Gene Variant Libraries: Design, Construction, and Research Applications

Rachel Speer, Ph.D.



# Gene Variant Libraries: Design, Construction, and Research Applications



- 1 Defining Gene Libraries and Mutant Libraries
- 2 Expression-Ready Gene Variant Libraries
- 3 Mutant Libraries: Rational, Systematic, and Random
- 4 Synthetic Biology Libraries

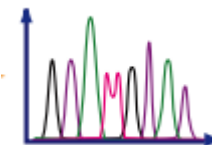
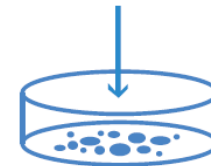
# GenScript – The most cited biology CRO



# Part 1: What are Gene Libraries?



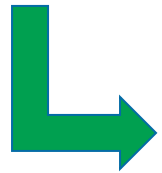
- a collection of many unique DNA sequences
- cloned into a vector
- propagated in micro-organisms
- screened for sequences of interest
- individually sequenced *before or after* assay



# The source of your insert DNA depends on your goal



**Inventory and study  
naturally-occurring  
DNA sequences**



Purify & clone nucleic acids  
from biological samples.  
*e.g. cDNA libraries,  
genomic libraries*

**Create  
novel sequences  
(synthetic genes)**



Construct desired  
sequences via  
mutagenesis, recombination,  
or gene synthesis.  
*e.g. Mutant Libraries*

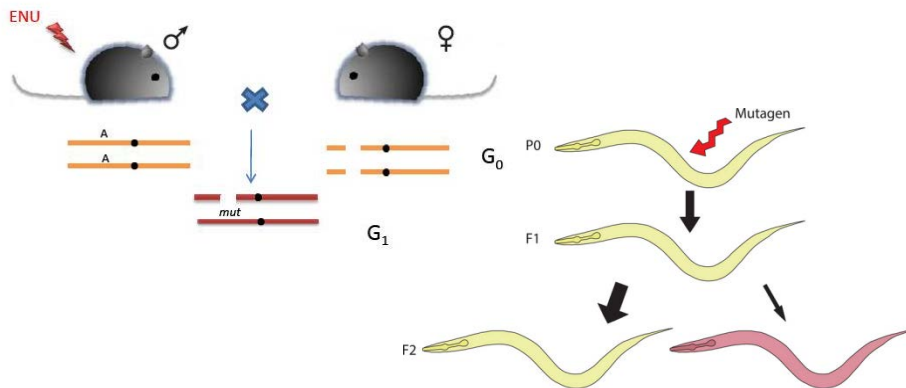
# Mutant Libraries



*in vivo* libraries

vs.

*in vitro* libraries



	50	60	70	80
	..... ..... ..... .....			
wildtype	GVVPILVELDGDVNGHKFSVSGEGEGDATYGKLTCLKFICT			
clone A01	GVVPILVELDGDVNGHKFSVSGEGEGDATYGKLTCLKFICT			
clone A02	GVVILVELDGDVNGHKFSVSGEGEGDATYGKLTCLKFICT			
clone A03	GVVPILVELDGDVNGHKFSVSGEGEGDATYGKLTCLKFICT			
clone A04	GVVPILVELDGDVNGHKFSVSGEGEGDATYGKLTCLKFICT			
clone A05	GVVPILVELDGDVNGHKFSVSGEGEGDATYGKLTCLKFICT			
clone A06	GVVPILVELDGDVNGHKFSVSGEGEGDATYGKLTCLKFICT			
clone A07	GVVPILVELDGDVNGHKFSVSGEGEGDATYGKLTCLKFICT			
clone A08	GVVPILVELDGDVNGHKFSVSGEGEGDATYGKLTCLKFICT			
....	....			
clone H12	GVVPILVELDGDVNGHKFSVSGEGEGDATYGKLTCLKFICT			

Expose organisms to a mutagen

Screen for phenotypes of interest

Identify gene mutation

Create mutant DNA sequences

Express in organism

Screen for desired phenotypes

# Synthetic Gene Libraries: Pools vs. Sequence-Verified Clones



## dsDNA Fragments

Mixture of sequences  
e.g. degenerated PCR products



## Transformants

(colonies = unique clones?)



## Sequence-verified clones

(purified plasmid DNA,  
glycerol stock)

- inexpensive & faster to create
- “quick & dirty” initial screening

**vs.**

- high-quality
- known identity
- reliable replication

# When to use synthetic gene libraries



- ◆ When you want to create something that doesn't already exist in nature.
- ◆ When you want to be systematic and unbiased
- ◆ When expressing genes in a different model organism
- ◆ When you know the sequences you want



# When not to use synthetic libraries



- ◆ When you want to create a library from a biological sample
- ◆ When you don't know what kind of variants you want



- ◆ Expression-Ready Gene Variant Libraries
- ◆ Mutant Libraries for Protein Engineering
- ◆ Synthetic Biology Libraries



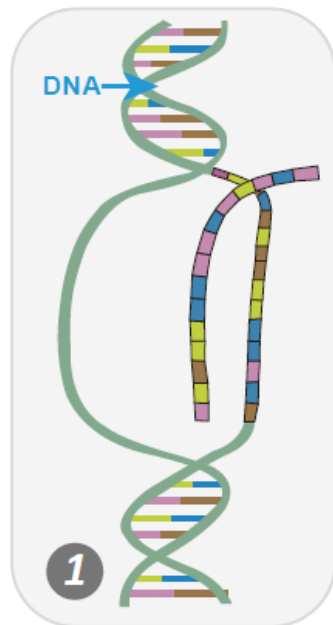
## ORFs cloned into expression vectors allow functional characterization of gene variants

- family members
- isoforms / splice variants
- disease-related variants
- mutants

## Why create synthetic libraries of naturally occurring genes?

- Ensure sequence accuracy
- Get expression-ready clones (optional: with tags)
- Improve expression levels through **codon optimization**

# Codon Optimization improves protein expression

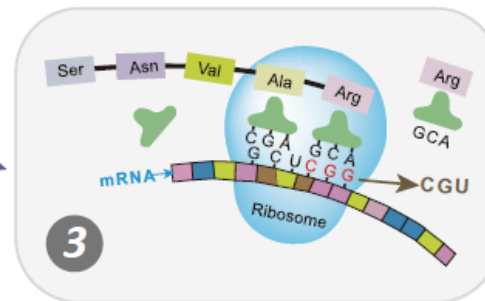


## 1. Transcription

- cis-regulatory elements (TATA box, termination signal, protein binding sites, etc.)
- chi sites
- polymerase slippage sites

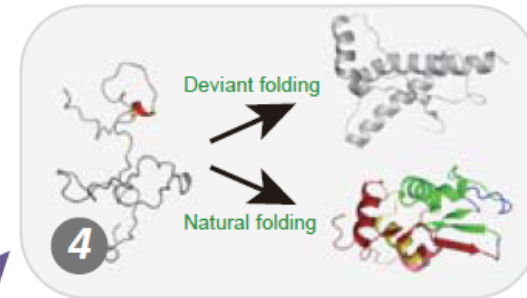
## 2. mRNA processing and stability

- cryptic splice sites
- mRNA secondary structure
- stable free energy of mRNA



## 3. Translation

- codon usage bias
- ribosomal binding sites (e.g. IRES)
- premature polyA sites



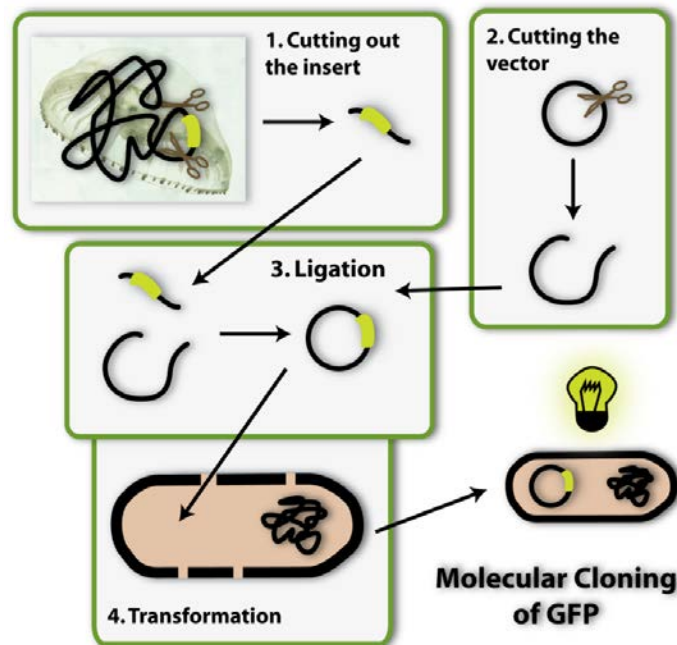
## 4. Protein folding

- codon context
- codon-anticodon interaction
- translation pause sites

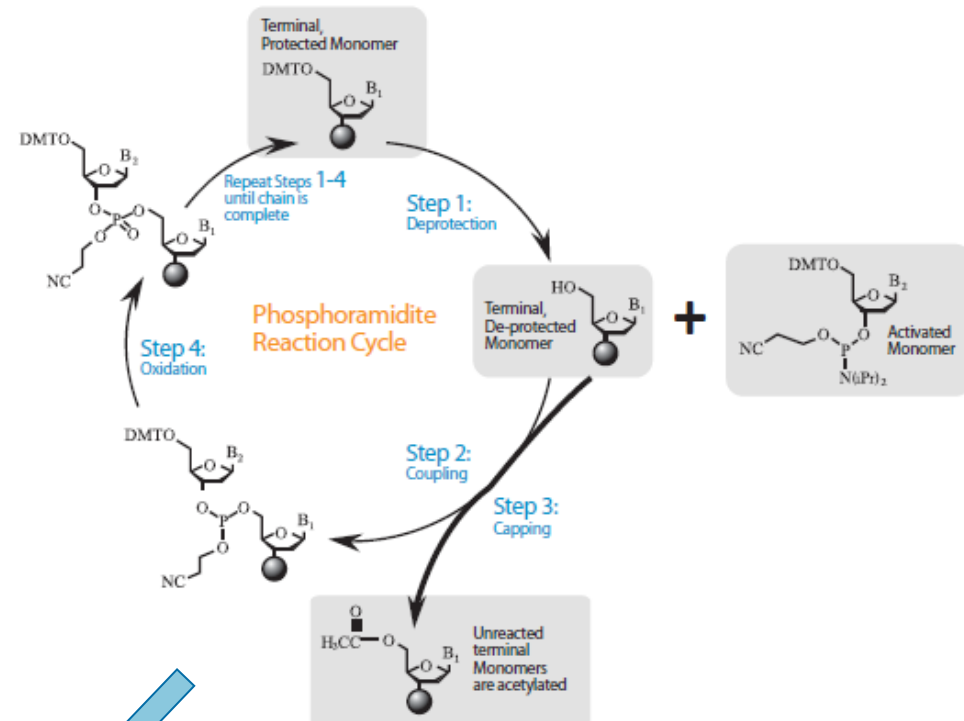
# Constructing Gene Libraries



Traditional Molecular Cloning:  
Cut & Paste naturally occurring sequences,  
PCR-based SDM as needed



Gene Synthesis:  
*de novo* Chemical Synthesis of DNA  
does not require a template



Clone into desired vector

# Choosing the most appropriate vector



Type of Cloning Vector	Advantages	Disadvantages
Plasmid	Replicate most prolifically; 700 copies per cell; most popular for inserts <5kb	Insert size Limited to 15kb (often lose larger inserts)
Bacteriophage	5–15kb	
Cosmids	replicate well, hold inserts 30–45kb, rarely lose inserts	Somewhat unstable, susceptible to recombination if contain repeats
BAC – Bacterial Artificial Chromosome	Insert size up to 350kb, fewer chimeras than YACs,	Only 1 copy per cell
YAC – Yeast Artificial Chromosome	Insert size up to 1000kb	High rate of chimeras

# Plasmids: cloning vectors vs. expression vectors



**Amplify, manipulate,  
and store DNA  
sequences**

**Efficiently express a  
sequence in my  
cell/tissue of interest**

## Cloning Vector

- high copy number
- multiple cloning region
- antibiotic resistance / lacZ for selection

## Expression Vector

- promoter for transcription
- Kozak (eukaryotic) or Shine-Dalgarno (prokaryotic) sequences for translation
- Viral vectors for *in vivo* delivery

You can create a library directly using your choice of Expression Vector; no need to shuttle!

# Plasmid-driven vs. endogenous expression



## Extra-chromosomal plasmid DNA



- strong promoter → overexpression
- *ideal for protein purification, reporter assays*

## Targeting Vector for integration into chromosomal DNA



- CRISPR/Cas9 is simpler and more efficient than ZFN, TALEN, Cre-lox
- *ideal to study gene function under endogenous promoter, normal stoichiometry*



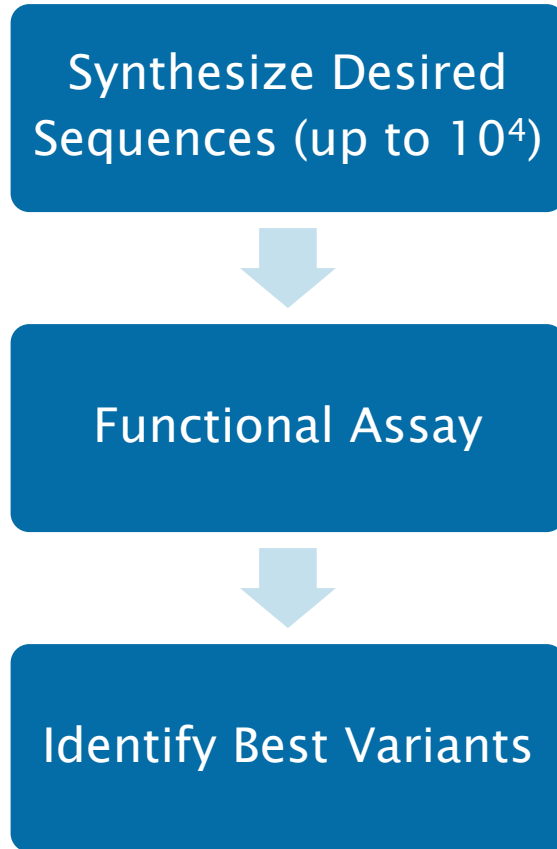


- ◆ Gene Variant Libraries
- ◆ Mutant Libraries
  - Rationally designed
  - Systematic / Saturated
  - Truncation variants
  - Random / degenerated
- ◆ Synthetic Biology Libraries

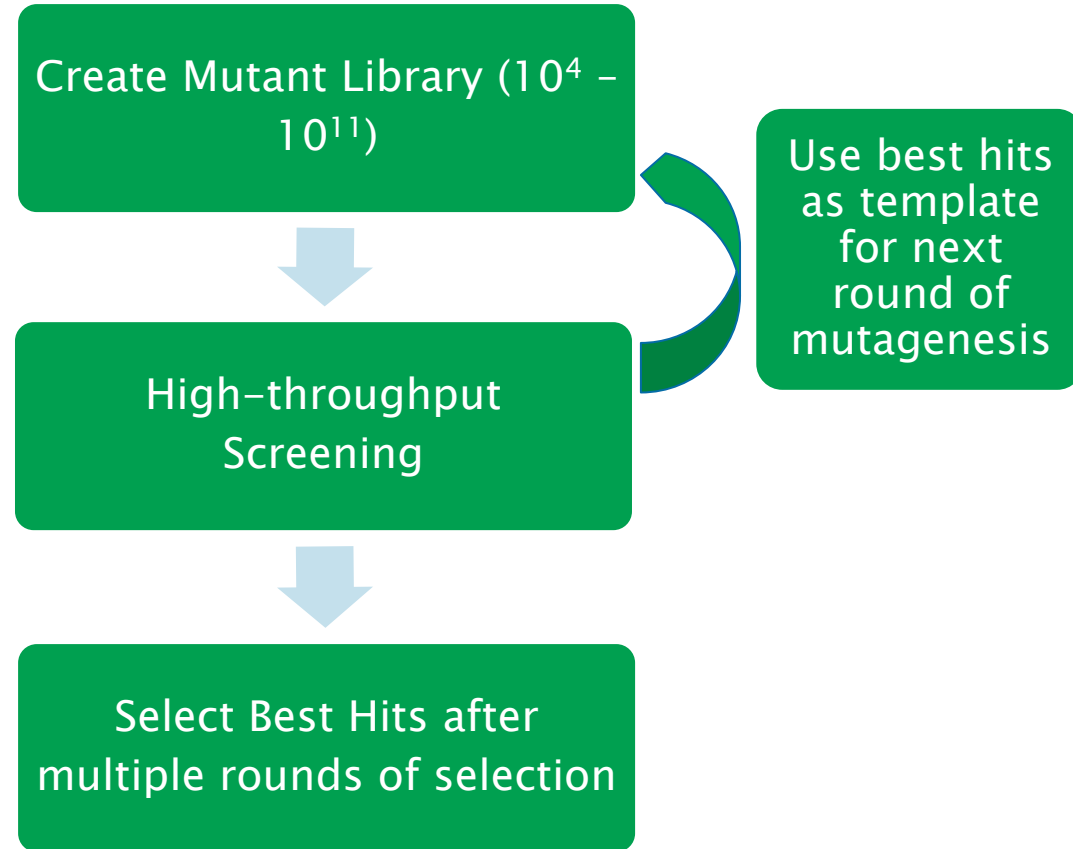
# Protein Engineering Strategies



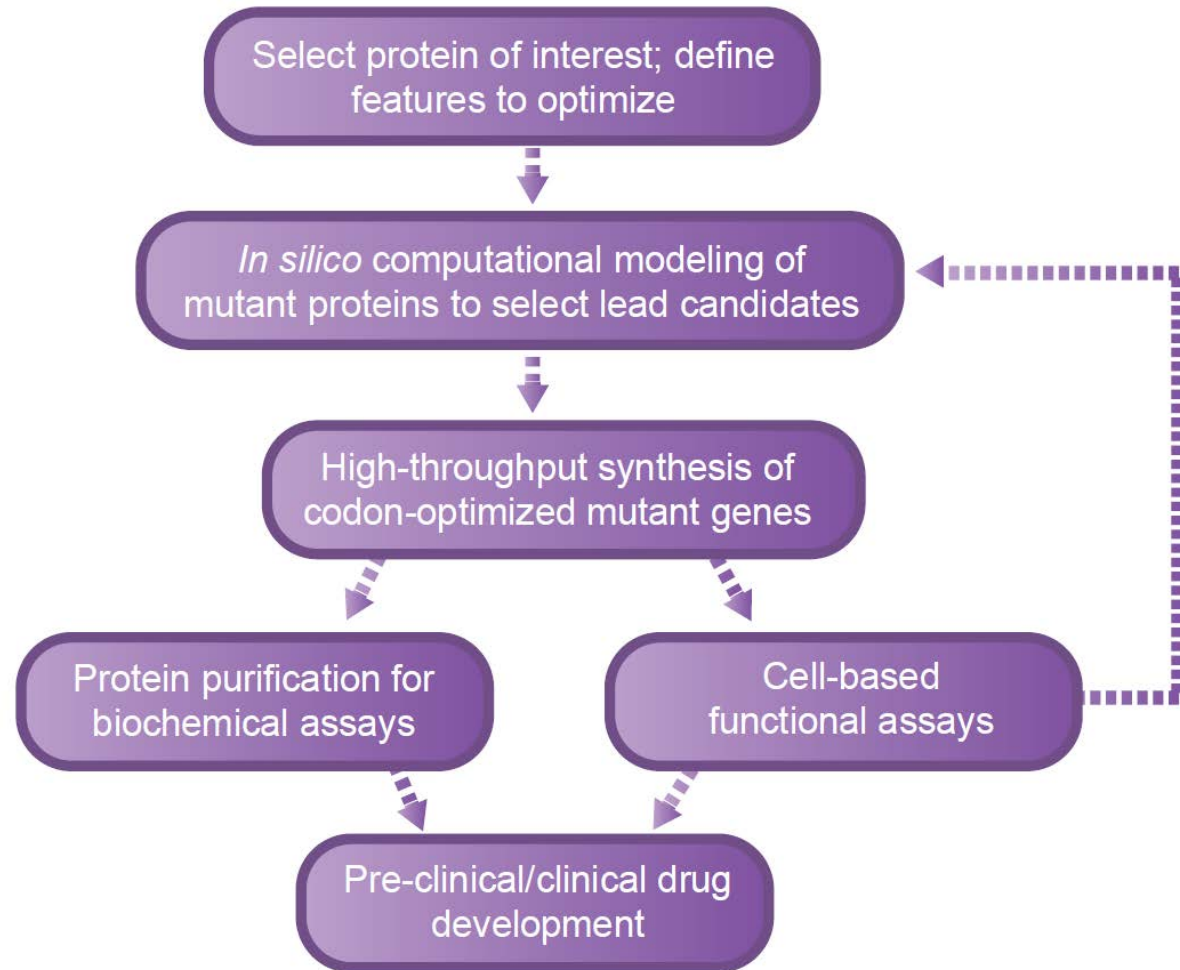
## Rational Design



## Directed Evolution



# Case Study: rational design of new drugs based on computational modeling



## Medical applications for protein engineering:

- vaccine design
- mAbs for immuno-oncology
- therapeutic enzymes

# Systematic Point Mutation Libraries



	50	60	70	80
<u>wildtype</u>	GVVPILVELDGDVNGHKFSVSGEGEDATYGKLTILKFICT			
F58A	GVVPILVELDGDVNGHKA SVSGEGEDATYGKLTILKFICT			
F58C	GVVPILVELDGDVNGHKC SVSGEGEDATYGKLTILKFICT			
F58D	GVVPILVELDGDVNGHKD SVSGEGEDATYGKLTILKFICT			
F58E	GVVPILVELDGDVNGHKE SVSGEGEDATYGKLTILKFICT			
F58G	GVVPILVELDGDVNGHKG SVSGEGEDATYGKLTILKFICT			
F58H	GVVPILVELDGDVNGHKH SVSGEGEDATYGKLTILKFICT			
F58I	GVVPILVELDGDVNGHKI SVSGEGEDATYGKLTILKFICT			
F58K	GVVPILVELDGDVNGHKK SVSGEGEDATYGKLTILKFICT			
....	....			
F58Y	GVVPILVELDGDVNGHKY SVSGEGEDATYGKLTILKFICT			

## Site-Saturation Mutagenesis

- a single amino acid
- multiple residues

```

1      M V S K G E E L F T G V V P I L X X L D G D V X G H
1      CCAIGGTGAGCAAAGGCGAAGAACTGTTTACCGGCGTGGTGCCGATTCTG NNSNNS CTGGATGGCGATGTG NNS GGCCAT
      GGTACCACTCGTTTCCGCTTCTTGACAAATGGCCGCACACCGGCTAAGACNNSNNSGACCTACCGCTACACNNSCCGGTA
      NcoI

81      K F S V S G E G E G D A T Y G K L T L K F I C T T G K
81      AAATTTAGCGTGAGCGGCGAAGGCGAAGGCGATGCGACCTATGGCAAACCTGACCCTGAAATTTATTTGCACCACCGGCAA
      TTTAAATCGCACTCGCCGCTTCCGCTTCCGCTACGCTGGATACCGTTTGACTGGGACTTTAAATAAACGTGGTGGCCGTT

161     L P V P W P T L V T T L T Y G V Q C F S R Y P D H M
161     ACTGCCGGTGCCGTGGCCGACCCCTGGTGACCACCCTGACCTATGGCGTGCAGTGCTTTAGCCGTTATCCGGATCATATGA
      TGACGGCCACGGCACCGGCTGGGACCACTGGTGGGACTGGATACCGCACGTACGAAATCGGCAATAGGCCTAGTATACT

241     K Q H D F F K S A M P X G Y V Q E R T I F F K D D G N
241     AACAGCATGATTTTTTTAAAGCGCGATGCGG NNSGGCTATGTGCAGGAACGTACCATTTTTTTTAAAGATGATGGCAAC
      TTGTCGTACTAAAAAATTTTCGCGCTACGGCNSCCGATACACGTCCTTGCAATGGTAAAAAATTTCTACTACCGTTG

321     Y K T R A X V K F E G D T L V N R X E L K G I D F K E
321     TATAAAACCGTGCC NNSGTGAATTTGAAGCGGATACCCCTGGTGAACCGT NNSGAACCTGAAAGGCATTGATTTTAAAGA
      ATATTTTGGGCACGCNNSCACTTTAAACTTCCGCTATGGGACCACTTGGCANNSTTGACTTTCCGTAATAAAATTTCT

401     D G N I L G H K L E Y N Y N S H N V Y I M A D K Q K
401     AGATGGCAACATTCTGGGCCATAAACTGGAATATACTATAACAGCCATAACGTGTATATTATGGCGGATAAACAGAAAA
  
```

# Scanning Libraries and Sequential Permutation Libraries



GVVPILVELDGDVNGHKFSVSGE GEGDATY GK  
GVVPILVELDGDVNG**A**KFSVSGE GEGDATY GK  
GVVPILVELDGDVNGH**A**FSVSGE GEGDATY GK  
GVVPILVELDGDVNGHK**A**SVSGE GEGDATY GK  
GVVPILVELDGDVNGHKF**A**VSGE GEGDATY GK

Alanine scanning

GVVPILVELDGDVNGHKFSVSGE GEGDATY GK  
GVVPILVELDGDVNG**X**KFSVSGE GEGDATY GK  
GVVPILVELDGDVNGH**X**FSVSGE GEGDATY GK  
GVVPILVELDGDVNGHK**X**SVSGE GEGDATY GK  
GVVPILVELDGDVNGHKF**X**VSGE GEGDATY GK

Consecutive  
site-saturations

GVVPILVELDGDVNGHKFSVSGE GEGDATY GK  
GVVPILVELDGDVNG**X**KFSVSGE GEGDATY GK  
GVVPILVELDGDVNG**XX**FSVSGE GEGDATY GK  
GVVPILVELDGDVNG**XXX**SVSGE GEGDATY GK  
GVVPILVELDGDVNG**XXXX**VSGE GEGDATY GK

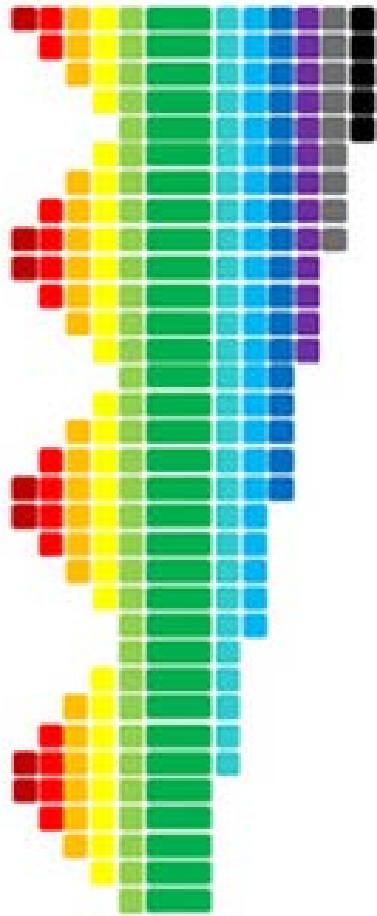
Sequential Permutation

# Size Matters: Constraining your design will reduce screening burden



Mutating only 6 residues of interest yields *many* unique sequences

- Site-saturated: 64 codons  $64^6 = 68,719,476,736$
- NNS:  $32^6 = 1,073,741,824$
- Trimer library: 20 a.a.  $20^6 = 64,000,000$
- Rationally designed SDM  $2*2*8*4*3*3 = 1152$



## Valuable for Structural Biology:

- Optimize protein solubility and stability
- Identify minimal domains required for folding, conformational stability, protein-protein interactions, catalytic activity

# Random Mutagenesis for Directed Evolution



	50	60	70	80
wildtype	GVVPILVELDGDVNGHKFSVSGEGEGDATYGKLTTLKFICT			
clone A01	GVVPILVELDGDVNGHKFSVSGEGEGDATYGKLTTLKFICT			
clone A02	GVVPILVELDGDVNGHKFSVSGEGEGDATYGKLTTLKFICT			
clone A03	GVVPILVELDGDVNGHKFSVSGEGEGDATYGKLTTLKFICT			
clone A04	GVVPILVELDGDVNGHKFSVSGEGEGDATYGKLTTLKFICT			
clone A05	GVVPILVELDGDVNGHKFSVSGEGEGDATYGKLTTLKFICT			
clone A06	GVVPILVELDGDVNGHKFSVSGEGEGDATYGKLTTLKFICT			
clone A07	GVVPILVELDGDVNGHKFSVSGEGEGDATYGKLTTLKFICT			
clone A08	GVVPILVELDGDVNGHKFSVSGEGEGDATYGKLTTLKFICT			
....	....			
clone H12	GVVPILVELDGDVNGHKFSVSGEGEGDATYGKLTTLKFICT			

Controlled randomization  
of complete reading frame

	50	60	70	80
wildtype	GVVPILVELDGDVNGHKFSVSGEGEGDATYGKLTTLKFICT			
clone A01	GVVPILVELDGDVNGHKFSVSGEGEGDATYGKLTTLKFICT			
clone A02	GVVPILVELDGDVNGHKFSVSGEGEGDATYGKLTTLKFICT			
clone A03	GVVPILVELDGDVNGHKFSVSGEGEGDATYGKLTTLKFICT			
clone A04	GVVPILVELDGDVNGHKFSVSGEGEGDATYGKLTTLKFICT			
clone A05	GVVPILVELDGDVNGHKFSVSGEGEGDATYGKLTTLKFICT			
clone A06	GVVPILVELDGDVNGHKFSVSGEGEGDATYGKLTTLKFICT			
clone A07	GVVPILVELDGDVNGHKFSVSGEGEGDATYGKLTTLKFICT			
clone A08	GVVPILVELDGDVNGHKFSVSGEGEGDATYGKLTTLKFICT			
....	....			
clone H12	GVVPILVELDGDVNGHKFSVSGEGEGDATYGKLTTLKFICT			

Controlled randomization  
of partial reading frame



# Techniques for constructing large (systematic or random) mutant libraries



## Error-prone PCR for random mutations

- Polymerase, cations, dNTPs, cycles

## Degenerate Oligos for site-saturation / sequential permutation libraries

- Single or combinations

## Trimers for efficient protein engineering

- Replace codons instead of single nucleotides

# Which library type is best?



Depends on:

- Your starting knowledge
- Your goals
- Your screening capacity

*Bigger isn't always better!*



# Rationally designed Mutant library for Promoter-Bashing



**Research Field:** Agriculture/Plant Biology

**Challenge:** Validate a new computational algorithm for identifying novel gene regulatory sequence motifs by systematic mutation

**Solution:** Synthesize a small library of constructs harboring systematic mutations in the 5' intergenic sequence.

GR2A										
1	2	3	4	5	6	7	8	9	10	11
GGCGCGCACC	ACCTGGGGCC	GGTACGTCGG	GAAGTSCCCA	CGCCTGGGCA	CGTCCAGCTT	TTCGTTGAAG	GATACCTGCG	TCTGCCACCT	ACGSCCACTA	CGGCGCGCGT
TTATATACAA	CAAGTTTTAA	TTGCATGATT	TCCTGTAAAC	ATAAGTTTAC	ATGAAGTAGG	GGATGGTCCT	TCGCAAGTAT	GAGTAACAAG	CATAAACAGC	ATTATATATG
12	13	14	15	16	17	18	19	20		
CCCGCCCGCG	TCTAAACAAG	CCCAAACGGC	CTTTAAACGG	CCTCGGTTCT	CGCAGTAGGS	GGGCCCTCTCT	TGACAGGGGG	AACACTGTCC	CAAAGCCTTC	CCAGGATTAC ...
AAATAAATAT	GAGGCCCACT	AAACCCATTA	AGGGCCCATI	AAGATTGGAG	ATACTGCTTT	TTTAAGAGAG	GTCACTTTTT	CCACAGTGAA	ACCCTAAGGA	AACTTCGGCA ...
								^TAIR	^EST	
GR11A										
		2	3	4	5	6	7	8	9	
		GGGGGCGGGG	AGGGGGGGGC	GCCAGGGGTT	GTCTAGGCCG	TTAAACCGCT	CAGTAGAGTG	TCCCTAAACC	CCAGCCTAAA	
		TTTTTATTTT	CTTTTTTTTA	TAACTTTTGG	TGAGCTTAAT	GGCCCAATAG	ACTGCTCTGT	GAAAGCCCAA	AACTAAGCCC	
10	11	12	13	14	15	16				
CCCGCCCGCG	CTTTGCTGCC	ATGCCGGTCT	AGCCTCCCAA	AGCTCTTGAG	AAGGATAAGC	ACCCCGAAAA	CGGGGTCGCC	GAGGACTACT	AATGGTAAGA	CCCCT
AAATAAAATA	AGGGTAGTAA	CGTAATTGAG	CTAAGAAACC	CTAGAGGTCT	CCTTCGCCTA	CAAAATCCCC	ATTTTGATAA	TCTTCAGCAG	CCGTTGCCTC	AAAAG
								^TAIR	^EST	

Davis IW, Benninger C, Benfey PN, and Elich T. POWRS: Position-Sensitive Motif Discovery PLoS One. 2012; 7(7): e40373. doi: 10.1371/journal.pone.0040373

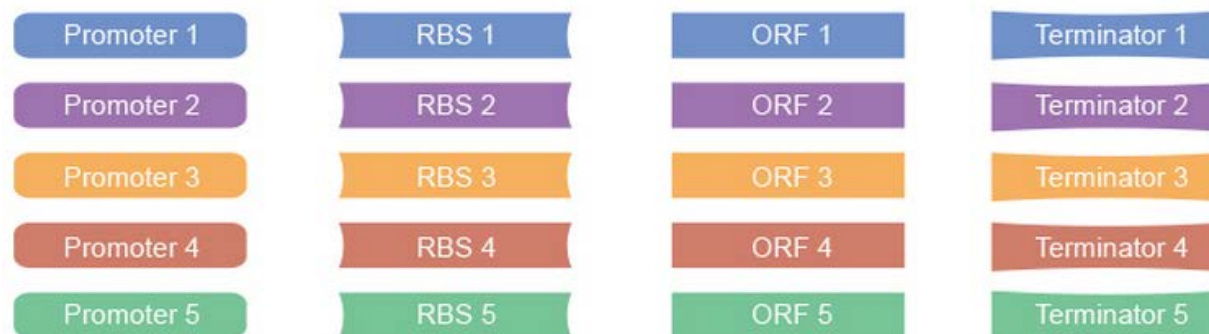


- ◆ Gene Variant Libraries
- ◆ Mutant Libraries
- ◆ Synthetic Biology Libraries
  - Combinatorial assembly variants
  - Synthetic Genomes / Synthetic Organisms

# Combinatorial Libraries

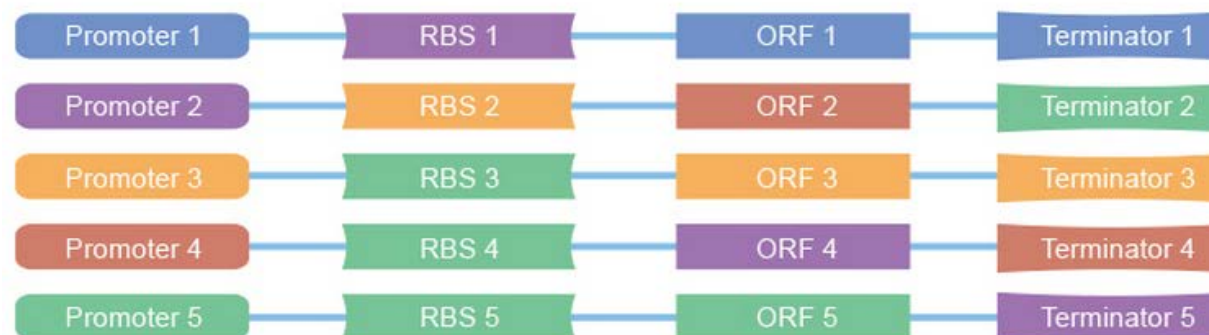


If you want to use 5 variants of each component:



Then you will need a combinatorial assembly library containing  $5 \times 5 \times 5 \times 5 = 625$  unique composite sequences

Just a few of those combinations:



Synthetic Biology  
based on standard parts

# Case Study: combinatorial assembly library for gene targeting / genome editing

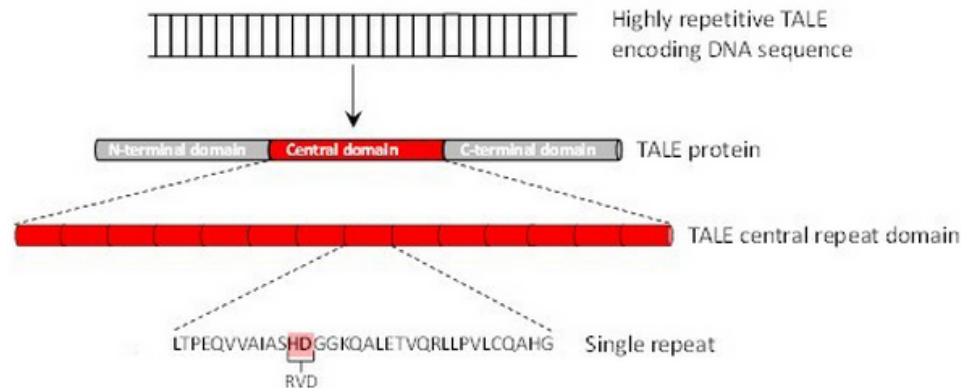


Kim, Y. *et al.* A library of TAL effector nucleases spanning the human genome. *Nat. Biotechnol.* **31**, 251–258 (2013).

- created TALENs for 18,740 unique protein-coding human genes
- synthesized 84 TALE plasmids containing all possible RVD combinations

**Table I. TALE specificity code**

RVD	Nucleotide
NI	adenine
HD	cytosine
NG	thymine
NN	guanine



**Fig 2. TALE protein organization**

# Codon-optimized gene synthesis solves problems for TALEN Combinatorial Library



1. Limit sequence similarity
2. Exclude rare codons to maximize translation efficiency
3. Guarantee accuracy of highly-repeated sequences



# Synthetic Genomes: a new model for *in vivo* mutant libraries



*GenScript is proud to be a contributing partner in the Sc2.0 International Consortium whose goal is to build a designer synthetic eukaryotic genome.*

## Sc2.0 – The Synthetic Yeast Genome Project

SCRaMbLE, a built-in inducible diversity generator

### Goals:

- Deduce the limits of chromosome structure
- Accelerate new strain development for biofuels & medical applications

### Building a Synthetic Eukaryotic Genome – Sc2.0



*Presented by: Leslie Mitchell, Ph.D., NYU Langone Medical Center*

June 25, 2014/  
2:00 pm EST

[Register now](#)



# Case Study: combining multiple design strategies in an “incognito library”



King, S. R. F. *et al.* Phytophthora infestans RXLR Effector PexRD2 Interacts with Host MAPKKK{varepsilon} to Suppress Plant Immune Signaling. *Plant Cell* **26**, 1345–1359 (2014).

## Design strategies:

- 1) **gene variants:** multiple PexRD2-like family members
- 2) **structure-led mutagenesis:** 5 specific PexRD2 mutants
- 3) **combinatorial assembly:** mutant genes fused with GFP, FLAG, or YN tags for different assays

## Construction method:

A combination of *de novo* gene synthesis, SDM, and recombination

# GenScript Toolkit



Library Type	Service	Advantages
Expression-Ready Clones	GenEZ ORF Cloning or Custom Cloning	<ul style="list-style-type: none"> <li>•choose your vector</li> <li>•10µg sequence-verified plasmid DNA</li> </ul>
Codon-Optimized Genes	OptimumGene + Gene Synthesis	<ul style="list-style-type: none"> <li>•increases protein expression</li> <li>•optional: protein expression evaluation</li> </ul>
Sequence-Verified Synthetic Library <i>25-10,000 variants</i>	GenPlus Next-Gen Gene Synthesis	<ul style="list-style-type: none"> <li>•cost-effective HT platform</li> <li>•4µg sequence-verified plasmid DNA</li> </ul>
Mutant Library – Rationally Designed <i>10<sup>1</sup>-10<sup>6</sup> variants</i>	Site-Directed Mutagenesis Library, Scanning Point Mutation	<ul style="list-style-type: none"> <li>•Choose sequence-verified clones or mixed library</li> </ul>
Mutant Library – Systematic or Randomized <i>10<sup>3</sup> - 10<sup>11</sup> variants</i>	Sequential Permutation Libraries or Randomized and Degenerated Libraries	<ul style="list-style-type: none"> <li>•10µg of dsDNA up to 10<sup>11</sup> variants</li> <li>•pooled clones in your choice of vector</li> <li>•pooled glycerol stock up to 10<sup>9</sup></li> </ul>
Truncation Variants	Truncation Variant Library	<ul style="list-style-type: none"> <li>•Up to 2000 variants within 4 weeks</li> <li>•4µg sequence-verified plasmid DNA</li> </ul>
Combinatorial Variants	Combinatorial Assembly	
Genome / Chromosome / metabolic circuits	GeneBricks	<ul style="list-style-type: none"> <li>•~10kb building blocks</li> <li>•100% guaranteed sequence accuracy</li> </ul>



- ◆ Thank you for attending!
- ◆ Please submit questions through chat.
- ◆ Please complete the survey you'll receive by email.
- ◆ Check upcoming and archived webinars at [www.genscript.com/webinars.html](http://www.genscript.com/webinars.html)
- ◆ Email me any time: [rachel.speer@genscript.com](mailto:rachel.speer@genscript.com)