

Experimental interrogation of the path dependence and stochasticity of protein evolution using phage-assisted continuous evolution

Bryan C. Dickinson, Aaron M. Leconte, Benjamin Allen, Kevin M. Esvelt, and David R. Liu

Supporting Information

Supporting Materials and Methods

General methods. All PCR reactions were performed with HotStartPhusion II polymerase (Thermo Scientific). Water was purified using a MilliQ water purification system (Millipore, Billerica MA). All vectors were constructed by isothermal assembly cloning (New England Biolabs). Single point mutants and reversions were constructed using the Quikchange II Site-Directed Mutagenesis kit (Agilent). All DNA cloning was performed with NEB Turbo cells (New England Biolabs). Plaque assays and PACE experiments were performed using *E. coli* S109 cells derived from DH10B as previously described (1). Luciferase assays were performed in NEB 10-beta cells (New England Biolabs).

Phage-Assisted Continuous Evolution (PACE). Host cell cultures, lagoons, media, and the PACE apparatus were as previously described (1). The accessory plasmids (APs) were constructed by modifying the previously described (1) low-copy AP with an SC101 origin to include a promoter of interest driving a tandem gIII-luciferase cassette, which were then also used for activity assays on isolated clones. To induce mutagenesis, 1% arabinose (wt/vol) was added to all populations throughout all PACE experiments (1). During the single-promoter stages of the evolution the lagoon volumes were fixed at 40 mL and during mixing stages of the evolutions the lagoon volumes were raised to 80 mL to keep the dilution rate constant. For the additional 24-h selection

on the final promoter, host cells flowed through each 40-mL lagoon at 3.5 volumes per hour. Lagoon samples were collected every 24 hours during the course of the evolution.

Phage pre-optimization. To minimize the potential fitness advantages of mutations to the phage genome, a previously described VCM13 helper phage with T7 RNAP (HP-T7RNAP) (1) was pre-optimized using PACE. Briefly, HP-T7RNAP was propagated for 6 days using a high-copy AP in which gene III expression is driven by a T7 promoter. Wild-type T7 RNAP was then subcloned into a randomly chosen phage backbone clone from the pre-optimization population and sequenced to ensure correct cloning of the T7 RNAP gene. All evolutions began with 40 μ L of 10^9 pfu/mL of this wild-type T7 RNAP phage.

Mutagenesis during PACE. The basal mutation rate of replicating filamentous phage (5.3×10^{-7} substitutions/bp) (2) is sufficient to generate all possible single but not double mutants of a given gene in a 100 mL lagoon following one generation of phage replication. For a target gene 1,000 base-pairs in length, a basal mutation rate of 5×10^{-7} applied to 5×10^{10} copies of the gene in a 100 mL lagoon yields 2.5×10^7 base substitutions, easily enough to cover all 3,000 single point mutants but not all double mutants. Arabinose induction of the MP can increase the mutation rate to $\sim 5 \times 10^{-5}$, yielding $\sim 2.5 \times 10^9$ mutations spread over 5×10^{10} copies of the gene after one generation. The vast majority of these are single mutants which together comprise a target area of $\sim 2.5 \times 10^{12}$ base pairs. Approximately 1.2×10^8 mutations should arise in a sequence of this length, sufficient to cover all 9×10^6 possible double mutants.

Plaque assays. S109 cells co-transformed with the MP and an AP of interest were grown in LB media to an OD_{600} of 0.8-1.0. 75 μ L of cells were then added to 25 μ L dilutions of phage pre-filtered with a 0.2 μ m syringe filter. After incubation for 5 min at 22 $^{\circ}$ C, 750 μ L of warm top agar (7 g/L agarose in LB) was added to the phage/cell

mixture, mixed by pipetting up and down twice, and plated onto quartered plates that had been previously poured with 4 mL of bottom agar (16 g/L agarose) in each quadrant. The plates were then grown overnight at 37 °C before plaques were counted.

T7 RNAP subcloning into expression vectors. Phage DNA was isolated from lagoon aliquots using a miniprep kit (Qiagen). The T7 RNAP genes were then amplified by PCR using primers that installed *Xho* I and *Bam*H I restriction sites. The previously described “expression plasmid” (EP) (1) was also amplified with primers that installed these restriction sites and, following PCR, was digested with *Dpn* I to remove remaining EP plasmid. The resulting PCR products were purified using a PCR purification kit (Qiagen), digested with *Xho* I and *Bam*H I, and purified by gel electrophoresis on a 0.7% agarose gel. 150 ng of each T7 RNAP library product was then ligated with 50 ng of each EP product using T4 DNA ligase. After PCR purification, the resulting libraries were transformed into S109 cells and plated on agar plates containing spectinomycin (50 µg/mL). After overnight growth, single colonies were picked from the plate, grown overnight, and processed using a miniprep kit to isolate EPs encoding individual T7 RNAP library members.

PCR primers for amplification of T7 RNAP genes:

5'-CAACATTCAAGGATCCACGGAATACCCAAAAGAACTGGCATG-3'

5'-AATCATCACACTCGAGCGGCGCAACTATCGGTATCAAGC-3'

PCR primers for amplification of EP:

5'-CAACATTCAACTCGAGATTTGAAGAGATAAATTGCACTGAAATCTAGAGCGG-3'

5'-AATCATCACAGGATCCAGAAAAGGAAGAGTATGAGGGAAGC-3'

Luciferase assays. EPs were co-transformed with an AP of interest into NEB 10-beta cells and plated onto agar plates (16 g/L in LB) with 50 µg/mL carbenicillin and 50

$\mu\text{g}/\text{mL}$ spectinomycin. After overnight growth at $37\text{ }^{\circ}\text{C}$, each well of a 96-well deep well plate containing 1 mL of LB with $50\text{ }\mu\text{g}/\text{mL}$ carbenicillin and $50\text{ }\mu\text{g}/\text{mL}$ spectinomycin was inoculated with a single colony. After growth with shaking at $37\text{ }^{\circ}\text{C}$ for 3 hours, $150\text{ }\mu\text{L}$ of each culture was transferred to a 96-well black wall, clear bottom plate (Costar). $1.2\text{ }\mu\text{L}$ of 10% decanal in ethanol (vol/vol) was then added to each well, the plate was shaken at room temperature for 1 minute, and then luminescence was measured on a Top Count NXT (Perkin Elmer) or Infinite M1000 Pro microplate reader (Tecan). The OD_{600} of each well was measured on a Spectramax M5 microplate reader (Molecular Devices) or an Infinite M1000 Pro microplate reader (Tecan).

The OD_{600} of a well containing only media was subtracted from all sample wells to obtain a corrected OD_{600} value for each well. The raw luminescence value for each well was then divided by that well's corrected OD_{600} value to obtain the luminescence value normalized to cell density. For the population level assays (Figs. 2A, 2B, 2C, S3), each single clone was analyzed in duplicate from two separately picked bacterial colonies and the resulting values were averaged to obtain a single point on the scatter plots. For all other assays, each sample was measured in triplicate, from three separate bacterial colonies, and the error bars shown are the standard error of those three independent measurements. For normalization, each plate contained triplicate samples of wild-type T7 RNAP co-transformed with T7-AP along with triplicate samples of an empty vector (without any T7 RNAP gene) co-transformed with the T7-AP and other APs of interest. The resulting average values for each sample were divided by the average value for the wild-type T7 RNAP acting on the T7-AP and multiplied by 100 to obtain the percent transcriptional activity relative to wild-type T7 RNAP on the T7 promoter. We validated that picking single clones and growing for 3 hours gives the same result as inoculating with an overnight culture (Fig. S16).

High-throughput sequencing (HTS) sample preparation. Lagoon aliquots were processed by Miniprep kit to isolate SP samples (Qiagen). The T7 RNAP genes, along with flanking DNA sequence both upstream and downstream, were amplified by PCR using the following primers:

5'-GGAGCAGGTCGCGGATTTTCG-3'

5'-GTCAAAAATGAAAATAGCAGCCTTTACAGAGAGAATAACATAAA-3'

The resulting PCR products were purified by gel electrophoresis on a 1% agarose gel and prepared for HTS using a Nextera kit (Illumina) and a slightly modified procedure. Briefly, 4 μL of DNA (2.5 $\text{ng}/\mu\text{L}$), 5 μL TD buffer, and 1 μL TDE1 were mixed together and then heated to 55 °C for 5 min. After purification (Zymo DNA purification kit), the resultant “tagmented” DNA samples were amplified with Illumina-supplied primers using the manufacturer’s protocol. The resulting PCR products were then purified using AMPure XP beads and the final concentration of DNA was quantified using PicoGreen (Invitrogen) and qPCR. The samples were sequenced on a MiSeq Sequencer (Illumina) in 2x150 paired-end runs using the manufacturer’s reagents following the manufacturer’s protocols.

High-throughput sequencing data analysis. A custom MATLAB script (available upon request) was used to align MiSeq reads to the wild-type sequence and count the nucleotide and amino acid positions from which the experimental sample deviates from the wild-type sequence. To compensate for systemic sample preparation and sequencing errors, the observed fraction of mutations at each nucleotide or amino acid position of the wild-type T7-RNAP reference gene was subtracted from the fraction of mutations in a given experimental sample to result in the “corrected fraction mutated” (3). Mutations were defined as amino acid positions with a corrected fraction mutation that is both $\geq 2.5\%$ and at least five standard deviations higher than the corrected fraction mutation of the wild-type reference sequence.

Calculation of Inverse Simpson Index (ISI). To quantify the diversity within each population, we used the inverse Simpson index (ISI) (4). At a single locus, this index is calculated as:

$${}^2D = \frac{1}{\sum_{\text{alleles } j} p_j^2}.$$

This formula can be extended to multiple loci by taking the harmonic mean over the single-locus values of 2D . The harmonic mean is more mathematically natural than the arithmetic mean for the ISI (5).

The ISI is an “effective number” (5,6). In our case, 2D is the effective number of alleles at a locus (in the single-locus version), or the average effective number of alleles across multiple loci (when the harmonic mean is taken over these loci). The number is “effective” because it also takes the evenness of allele frequency distributions into account. For example, consider a locus with two alleles of frequency 99% and 1%. Though the actual number of alleles is two, the ISI gives an effective number of alleles of only ${}^2D \approx 1.02$, since one allele comprises the vast majority of the frequency. At the other extreme, if there are m equally abundant alleles, we have ${}^2D = m$.

In the case of our experiment, the ISI takes values between 1 and 2, since we consider only two alleles (wild-type and mutant) at each locus.

To obtain an unbiased estimate of the ISI from high-throughput sequencing measurements, we first compute an “effective sample size” \tilde{N} , as $\tilde{N} = 1/(4\epsilon^2)$, where ϵ is the average error per measurement per locus. With this effective sample size, the standard deviation due to sampling error in estimating an allele frequency is $2\epsilon\sqrt{p(1-p)}$, where p is the true frequency; this quantity has a maximum value of ϵ at $p = 1/2$. We then use Nei and Chesser’s (7) unbiased estimator for the convention Simpson index:

$$\hat{H} = \frac{2\tilde{N}}{2\tilde{N} - 1} \left(1 - \sum_{\text{alleles } j} p_j^2 \right),$$

where p_i is the measured frequency of allele i . Finally, we transform this into an unbiased estimator of the ISI:

$${}^2\hat{D} = \frac{1}{1 - \hat{H}}.$$

For a collection of populations (in our case, populations), α -diversity is the average diversity across populations. Specifically, we take the harmonic mean of the values of 2D across each population. The harmonic mean has better mathematical properties than the arithmetic mean for the ISI (5).

Calculation of F_{ST} . For a single locus with two alleles (wild-type and mutant), F_{ST} (8) is defined by the formula

$$F_{ST} = \frac{\text{Var}(p)}{\bar{p}(1 - \bar{p})},$$

where \bar{p} is the mean frequency of the mutant allele across subpopulations, and $\text{Var}(p)$ is the variance of this frequency across subpopulations. F_{ST} takes values between 0 and 1, with 0 indicating that all populations are identical in allele frequency. Large values of F_{ST} indicate greater divergence between subpopulations. We caution, however, that most theoretical work on F_{ST} concerns the case of weak selection, nonzero migration, and many generations (9,10). More modeling work is needed to understand the behavior of F_{ST} under strong selection, zero migration, and relatively few generations, as is the case for this experiment.

We calculated F_{ST} using the θ estimator of Weir and Cockerham (11). In referring to θ as an estimator of F_{ST} , we adopt Weir and Cockerham's view that the n subpopulations under consideration can be regarded as a sample from a much larger collection of subpopulations. In our case, this means that the population populations we observed are a subset of the infinite number of populations that could theoretically be created.

The θ estimator applies to the situation of estimating F_{ST} from a sample of size N taken from a large population. To relate this situation to ours, we calculated an "effective sample size" that would produce a sampling error of the same magnitude as the measurement error in the high-throughput sequencer. Specifically, we calculated the effective sample size as $\tilde{N} = 1/(4\epsilon^2)$, where ϵ is the average error per measurement per locus. With this effective sample size, the standard deviation due to sampling error

in estimating an allele frequency is $2\epsilon\sqrt{p(1-p)}$, where p is the true frequency; this quantity has a maximum value of ϵ at $p = 1/2$.

To compute θ , suppose that, at a specific locus, the frequencies of the mutant allele in the n subpopulations are measured as p_1, \dots, p_n . Then θ at this locus is given by (11):

$$\theta = \frac{s^2 - \frac{1}{\bar{N}-1} \left(\bar{p}(1-\bar{p}) - \frac{n}{n-1} s^2 \right)}{\bar{p}(1-\bar{p}) + \frac{s^2}{n}}.$$

Above, \bar{p} and s^2 are, respectively, the sample mean and sample variance of the frequency of the mutant allele across subpopulations:

$$\bar{p} = \frac{1}{n} \sum_{i=1}^n p_i,$$

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (p_i - \bar{p})^2.$$

In our calculations, each position in the protein sequence was regarded as a single locus. We obtained a multilocus F_{ST} value by summing the numerator and denominator in the formula for θ across loci before dividing, as suggested (11).

We estimated the error of F_{ST} due to measurement using the standard formula

$$\sigma = \epsilon \sqrt{\sum_{i=1}^n \left(\frac{\partial f}{\partial p_i} \right)^2}.$$

Above, σ is the standard deviation due to measurement error in f , and ϵ is the estimated standard deviation due to error in the measurement of each allele frequency p_i .

Phylogenetic analysis. Several clones from each lagoon were selected at random for Sanger sequencing at 96 h and 192 h (Figs. S4-S7). A phylogenetic tree was constructed from the amino acid sequences at each time point using Bayesian analysis (MrBayes 3.2.1) (12). The evolutionary model specified a single rate across sites, with

gamma-distributed rate variation with a proportion of invariant sites. Satisfactory convergence was assumed if the average standard deviation of split frequencies was <0.02 . At least four independent analyses were run and compared for each data set to confirm the tree topology for clades assigned with high confidence. Consensus trees were visualized in FigTree 1.4.

Supporting Discussion

Population genotypic characterization. To statistically evaluate the total genetic divergence between the evolved populations, we used F_{ST} , a widely applied measure of population differentiation that estimates the proportion of variation between populations that is beyond the variation seen within populations (13). Compared populations are more related as F_{ST} approaches 0 and are more divergent as F_{ST} approaches 1. We compared F_{ST} values for sibling populations (those subjected to identical selection conditions) with the F_{ST} value for all populations taken together, which represents the total genetic diversity of the system at a given time point. At 96 h, the T3 and SP6 populations exhibited significantly lower F_{ST} than the total F_{ST} for all eight populations (Fig. S10C), indicating that the populations have divergent, pathway-dependent genotypes. From 96 to 144 h the SP6 populations diverged, while the T3 populations converged, indicating that the T3 populations experienced a decrease in divergence as they adapted to the T3/final promoter stepping stone. Instead of converging once the full final promoter selection began at 144 h, the T3 populations diverged again, exhibiting internal divergence by 216 h comparable to the total divergence of all populations, which remained nearly constant throughout the convergent part of the selection. These results reveal that comparisons of population-wide genotypes are unable to account for the different phenotypic outcomes of the two pathways, and suggest that either mutational stochasticity or convergence was a primary determinant of population-wide measures of genetic evolution in this system.

Bayesian phylogenetic analysis on the complete single-clone data (Figs. S11 and S12) confirmed that variants from the same population tend to be genetically similar and

could often be identified as belonging to the same clade. However, variants from different populations following the same pathway (SP6 or T3) could not be confidently grouped with each other at 96 h or 192 h. Differences between population replicates confounded the differences between pathways, further indicating that mutational stochasticity between replicates is a primary determinant of genetic outcome.

SUPPORTING REFERENCES

1. Esvelt KM, Carlson JC, & Liu DR (2011) A system for the continuous directed evolution of biomolecules. *Nature* 472(7344):499-503.
2. Drake JW (1991) A constant rate of spontaneous mutation in DNA-based microbes. *Proc Natl Acad Sci U S A* 88(16):7160-7164.
3. Leconte AM, et al. (2013) A population-based experimental model for protein evolution: Effects of mutation rate and selection stringency on evolutionary outcomes. *Biochemistry* 52(8):1490–1499.
4. Simpson EH (1949) Measurement of Diversity. *Nature* 163:699.
5. Jost L (2007) Partitioning diversity into independent alpha and beta components. *Ecology* 88(10):2427-2439.
6. Jost L (2006) Entropy and diversity. *Oikos* 113(2):363-375.
7. Nei M & Chesser RK (1983) Estimation of fixation indices and gene diversities. *Ann Hum Genet* 47(Pt 3):253-259.
8. Wright S (1943) Isolation by Distance. *Genetics* 28(2):114-138.
9. Rousset F (2004) *Genetic Structure and Selection in Subdivided Populations* (Princeton University Press).
10. Whitlock MC (2011) G'ST and D do not replace FST. *Molecular Ecology* 20(6):1083-1091.
11. Weir BS & Cockerham CC (1984) Estimating F-Statistics for the Analysis of Population Structure. *Evolution* 38(6):1358-1370.
12. Ronquist F, et al. (2012) MrBayes 3.2: efficient Bayesian phylogenetic inference and model choice across a large model space. *Syst Biol* 61(3):539-542.
13. Holsinger KE, Weir BS (2009) Genetics in geographically structured populations: Defining, estimating and interpreting F(ST). *Nat Rev Genet* 10(9):639–650.

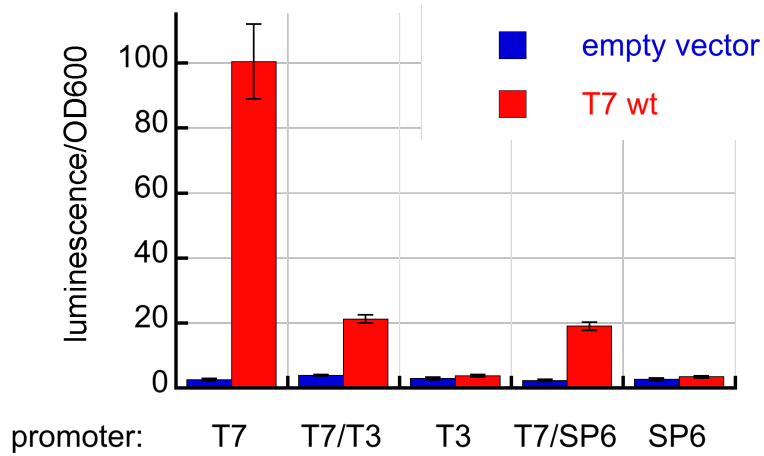


Fig. S1. Activity of wild-type T7 RNAP on intermediate and target promoters. Normalized luminescence of an empty vector control (blue bars) and wild-type T7 RNAP (red bars) on reporter vectors driven by the T7, T7/T3 hybrid, T3, T7/SP6 hybrid, and SP6 promoter. Error bars represent standard error (n = 4).

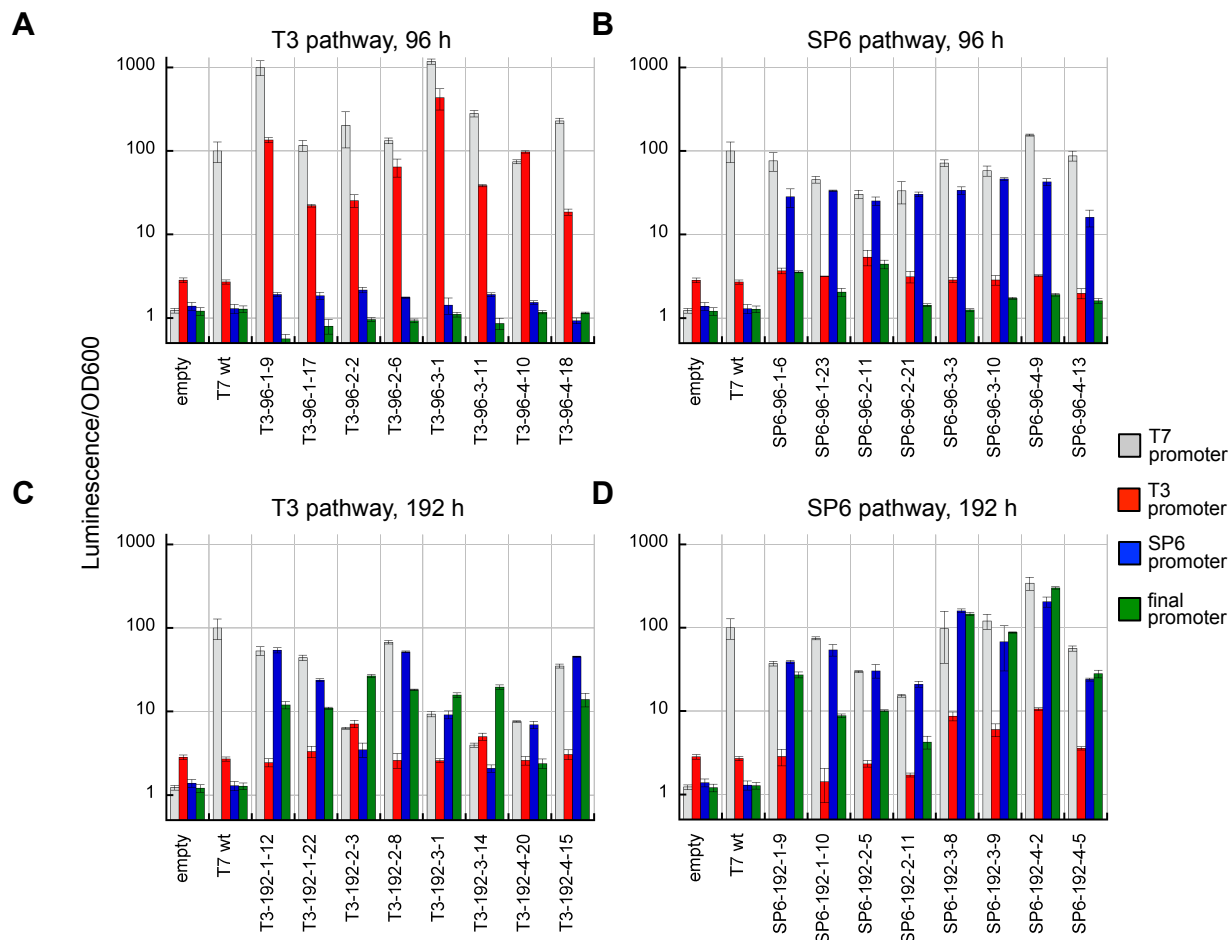


Fig. S2. Activity profiles of single variants. (A) Activity assays of two clones from each population from the 96 h time point of the T3 pathway. (B) Activity assays of two clones from each population from the 96 h time point of the SP6 pathway. (C) Activity assays of two clones from each population from the 192 h time point of the T3 pathway. (D) Activity assays of two clones from each population from the 192 h time point of the SP6 pathway. All data is normalized to wild-type T7 RNAP acting on the T7 promoter (100% by definition). Single clones correspond to the full sequences in Figs. S4 through S7. The nomenclature for clone ID is “pathway-time point-population #-variant #”. For example T3-96-1-9 is T3 pathway, 96 h time point, population 1, variant 9.

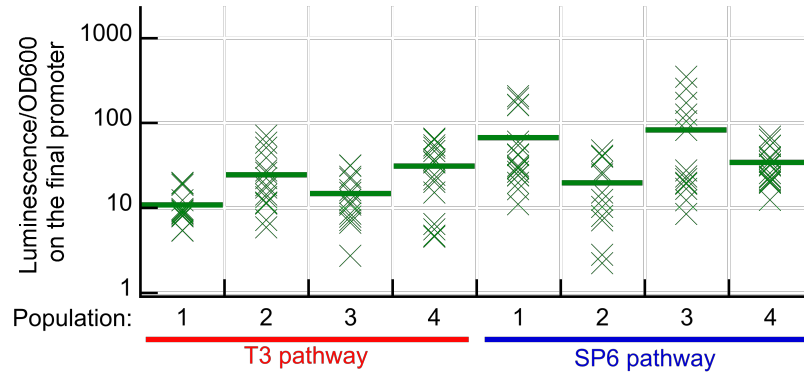


Fig. S3. Phenotypes of evolved populations at 216h. Normalized luminescence at 216 h of the T3 and SP6 populations on the final reporter vector. Each point represents the activity of a single randomly isolated clone.

evolved RNAP clone:					SP6-96-2-2	SP6-96-2-11	SP6-96-2-21	SP6-96-2-23	SP6-96-2-24	SP6-96-3-2	SP6-96-3-3	SP6-96-3-9	SP6-96-3-10	SP6-96-3-15	SP6-96-4-1	SP6-96-4-9	SP6-96-4-12	SP6-96-4-13	SP6-96-4-21
															N9H			N9H	
		A65T			F55L										M54V	M54V	M54V		M54V
Q107K	Q107K	Q107K	Q107K			A94S									D66N	D66N	D66N		D66N
			A138T				A124T												
				K163R															
		E218G			V174A														V214E
E222A	E222A		E222A	E222K	M219R	M219R	E222K	E222K	E222K	E222K	E222K	E222K	E222K	E222K	E222K	E222K	E222K	E222K	E222K
		V426F	K450R																
I479T	I479T		E504K	I479T															
S539F	S539F			S539F	S539F	S539F		S539F							C510R	C510R	C510R		C510R
		M635I													A575V	A575V	A575V		A575V
E643K	E643K		E643K	E643K															
V685A	V685A	V685A	V685A	T688M	V685A	V685A	V685A	V685A	T688M	F646L	F646L	F646L	F646L	F646L	V685A	V685A	V685A	V685A	V685A
										V685A	V685A	Q672K	V685A	V685A					
Q758K	Q758K	Q758K	Q758K	Q758K	Q758K	Q758K	Q758K	Q758K	Q758K	Q758K	Q758K	Q758K	Q758K	Q758K	A724V	A724V	A724V	A724V	A724V
A866T															Q758R	Q758R	Q758R	Q758K	Q758R

Fig. S5. Fully sequenced genes of clones from the 96-h time point of the SP6 pathway. The clone ID is shown at the top of each column and the mutations present in that clone listed down the column. The nomenclature for clone ID is “pathway-time point-population #-variant #”. For example SP6-96-1-2 is SP6 pathway, 96 h time point, population 1, variant 2.

evolved RNAP clone:					T3-192-1-1	T3-192-1-2	T3-192-1-12	T3-192-1-21	T3-192-1-22	T3-192-2-2	T3-192-2-3	T3-192-2-8	T3-192-2-9	T3-192-2-16	T3-192-3-1	T3-192-3-2	T3-192-3-5	T3-192-3-11	T3-192-3-14	T3-192-4-2	T3-192-4-3	T3-192-4-4	T3-192-4-5	T3-192-4-15	T3-192-4-20	
						N5T	N5T	N5T		E90K K93T		N86S	E90K								A25T E56D		A25T E56D	A25T E56D	A25T E56D	
					S128R				S128R V134I		D130N				S128N							T127I				
									M219K		N165S	H176R		N165S								D130N	D130N	D130N	D130N	D130N
					E222K	E222K	E222K	E222K		E222K	E222K			E222K	E222K	E222K	E222K	E222K	E222K	E222K	E222K	E222K	E222K	E222K	E222K	E222K
															Q239K		Q238S	Q232R								
					T256I	T256I	T256I					A247S				T243I	T243I		T243I							
					G263D	G263D	G263D			S301K			S301K	S301K												
										V334I		A373D														
					E403D	S397R	S397R		E403D			M439V				M431T										
																D537Y										
					G542V	G542V	G542V	G542V	G542V		G555C	G555C			G542V	G542V	G542V	G542V	G542V		G542V	G542V	G542V	G542V		G542V
										V574A			V574A	V574A		S564G		S564G							V574A	
											E643K			E600K		E643K	E643K	E643K	E643K		E643K	E643K	E643K	E643K		E643K
												Q656H	Q648L			G675R	Q656H	Q656H	G675R	Q656H						
												W682R				E883D		E883D								
										S686I				S686I		S686N	S686N		S686N							
					R756C Q758R	R756C Q758R	R756C Q758K	R756C Q758R	R756C Q758K D770N	N748D	N748D	R756C Q758K	N748D	N748D	R756C Q758R	N748D	N748D	N748D	R756C Q758R	N748D	N748D	R756C Q758K	R756C Q758K	R756C Q758K	R756C Q758K	R756C Q758K
					H772R	H772R	H772R	H772R		H772R			H772R	H772R												
											E775A V783I				E775V	E775V										
									L864F																	

Fig. S6. Fully sequenced genes of clones from the 192-h time point of the T3 pathway. The clone ID is shown at the top of each column and the mutations present in that clone listed down the column. The nomenclature for clone ID is “pathway-time point-population #-variant #”. For example T3-192-1-2 is T3 pathway, 192 h time point, population 1, variant 2.

gIII promoter Start

TTTAAGAAATTCACCTCGAAAGCAAGCTGATAAACCGATAACAATTAAGGCTCCTTTTGGAGCCTTTTTTCGCGCCAGAAGGAAACCATG

T3-96-1-9 TTTAAGAAATTCACCTCGAAAGCAAGCTGATAAACCGATAACAATTAAGGCTCCTTTTGGAGCCTTTTTTCGCGCCAGAAGGAAACCATG

T3-96-1-17 TTTAAGAAATTCACCTCGAAAGCAAGCTGATAAACCGATAACAATTAAGGCTCCTTTTGGAGCCTTTTTTCGCGCCAGAAGGAAACCATG

T3-96-2-2 TTTAAGAAATTCACCTCGAAAGCAAGCTGATAAACCGATAACAATTAAGGCTCCTTTTGGAGCCTTTTTTCGCGCCAGAAGGAAACCATG

T3-96-2-6 TTTAAGAAATTCACCTCGAAAGCAAGCTGATAAACCGATAACAATTAAGGCTCCTTTTGGAGCCTTTTTTCGCGCCAGAAGGAAACCATG

T3-96-3-1 TTTAAGAAATTCACCTCGAAAGCAAGCTGATAAACCGATAACAATTAAGGCTCCTTTTGGAGCCTTTTTTCGCGCCAGAAGGAAACCATG

T3-96-3-11 TTTAAGAAATTCACCTCGAAAGCAAGCTGATAAACCGATAACAATTAAGGCTCCTTTTGGAGCCTTTTTTCGCGCCAGAAGGAAACCATG

T3-96-4-10 TTTAAGAAATTCACCTCGAAAGCAAGCTGATAAACCGATAACAATTAAGGCTCCTTTTGGAGCCTTTTTTCGCGCCAGAAGGAAACCATG

T3-96-4-18 TTTAAGAAATTCACCTCGAAAGCAAGCTGATAAACCGATAACAATTAAGGCTCCTTTTGGAGCCTTTTTTCGCGCCAGAAGGAAACCATG

SP6-96-1-6 TTTAAGAAATTCACCTCGAAAGCAAGCTGATAAACCGATAACAATTAAGGCTCCTTTTGGAGCCTTTTTTCGCGCCAGAAGGAAACCATG

SP6-96-1-23 TTTAAGAAATTCACCTCGAAAGCAAGCTGATAAACCGATAACAATTAAGGCTCCTTTTGGAGCCTTTTTTCGCGCCAGAAGGAAACCATG

SP6-96-2-11 TTTAAGAAATTCACCTCGAAAGCAAGCTGATAAACCGATAACAATTAAGGCTCCTTTTGGAGCCTTTTTTCGCGCCAGAAGGAAACCATG

SP6-96-2-12 TTTAAGAAATTCACCTCGAAAGCAAGCTGATAAACCGATAACAATTAAGGCTCCTTTTGGAGCCTTTTTTCGCGCCAGAAGGAAACCATG

SP6-96-3-3 TTTAAGAAATTCACCTCGAAAGCAAGCTGATAAACCGATAACAATTAAGGCTCCTTTTGGAGCCTTTTTTCGCGCCAGAAGGAAACCATG

SP6-96-3-10 TTTAAGAAATTCACCTCGAAAGCAAGCTGATAAACCGATAACAATTAAGGCTCCTTTTGGAGCCTTTTTTCGCGCCAGAAGGAAACCATG

SP6-96-4-9 TTTAAGAAATTCACCTCGAAAGCAAGCTGATAAACCGATAACAATTAAGGCTCCTTTTGGAGCCTTTTTTCGCGCCAGAAGGAAACCATG

SP6-96-4-13 TTTAAGAAATTCACCTCGAAAGCAAGCTGATAAACCGATAACAATTAAGGCTCCTTTTGGAGCCTTTTTTCGCGCCAGAAGGAAACCATG

T3-192-1-12 TTTAAGAAATTCACCTCGAAAGCAAGCTGATAAACCGATAACAATTAAGGCTCCTTTTGGAGCCTTTTTTCGCGCCAGAAGGAAACCATG

T3-192-1-22 TTTAAGAAATTCACCTCGAAAGCAAGCTGATAAACCGATAACAATTAAGGCTCCTTTTGGAGCCTTTTTTCGCGCCAGAAGGAAACCATG

T3-192-2-3 TTTAAGAAATTCACCTCGAAAGCAAGCTGATAAACCGATAACAATTAAGGCTCCTTTTGGAGCCTTTTTTCGCGCCAGAAGGAAACCATG

T3-192-2-8 TTTAAGAAATTCACCTCGAAAGCAAGCTGATAAACCGATAACAATTAAGGCTCCTTTTGGAGCCTTTTTTCGCGCCAGAAGGAAACCATG

T3-192-3-1 TTTAAGAAATTCACCTCGAAAGCAAGCTGATAAACCGATAACAATTAAGGCTCCTTTTGGAGCCTTTTTTCGCGCCAGAAGGAAACCATG

T3-192-3-14 TTTAAGAAATTCACCTCGAAAGCAAGCTGATAAACCGATAACAATTAAGGCTCCTTTTGGAGCCTTTTTTCGCGCCAGAAGGAAACCATG

T3-192-4-15 TTTAAGAAATTCACCTCGAAAGCAAGCTGATAAACCGATAACAATTAAGGCTCCTTTTGGAGCCTTTTTTCGCGCCAGAAGGAAACCATG

T3-192-4-20 TTTAAGAAATTCACCTCGAAAGCAAGCTGATAAACCGATAACAATTAAGGCTCCTTTTGGAGCCTTTTTTCGCGCCAGAAGGAAACCATG

SP6-192-1-9 TTTAAGAAATTCACCTCGAAAGCAAGCTGATAAACCGATAACAATTAAGGCTCCTTTTGGAGCCTTTTTTCGCGCCAGAAGGAAACCATG

SP6-192-1-10 TTTAAGAAATTCACCTCGAAAGCAAGCTGATAAACCGATAACAATTAAGGCTCCTTTTGGAGCCTTTTTTCGCGCCAGAAGGAAACCATG

SP6-192-2-1 TTTAAGAAATTCACCTCGAAAGCAAGCTGATAAACCGATAACAATTAAGGCTCCTTTTGGAGCCTTTTTTCGCGCCAGAAGGAAACCATG

SP6-192-2-11 TTTAAGAAATTCACCTCGAAAGCAAGCTGATAAACCGATAACAATTAAGGCTCCTTTTGGAGCCTTTTTTCGCGCCAGAAGGAAACCATG

SP6-192-3-8 TTTAAGAAATTCACCTCGAAAGCAAGCTGATAAACCGATAACAATTAAGGCTCCTTTTGGAGCCTTTTTTCGCGCCAGAAGGAAACCATG

SP6-192-3-9 TTTAAGAAATTCACCTCGAAAGCAAGCTGATAAACCGATAACAATTAAGGCTCCTTTTGGAGCCTTTTTTCGCGCCAGAAGGAAACCATG

SP6-192-4-2 TTTAAGAAATTCACCTCGAAAGCAAGCTGATAAACCGATAACAATTAAGGCTCCTTTTGGAGCCTTTTTTCGCGCCAGAAGGAAACCATG

SP6-192-4-5 TTTAAGAAATTCACCTCGAAAGCAAGCTGATAAACCGATAACAATTAAGGCTCCTTTTGGAGCCTTTTTTCGCGCCAGAAGGAAACCATG

Fig. S8. Sequencing data for promoters and RBS of single variants. The upstream DNA sequences containing the promoter and RBS driving the expression of the subset of clones assayed were fully sequenced. Although several mutations arise in different clones, they are generally not located in the promoter or RBS and do not correlate with the activity differences of clones. The phage genome, including the promoter and RBS driving T7 RNAP activity, were pre-evolved for 6 days, during which time these sequences had the opportunity to be optimized.

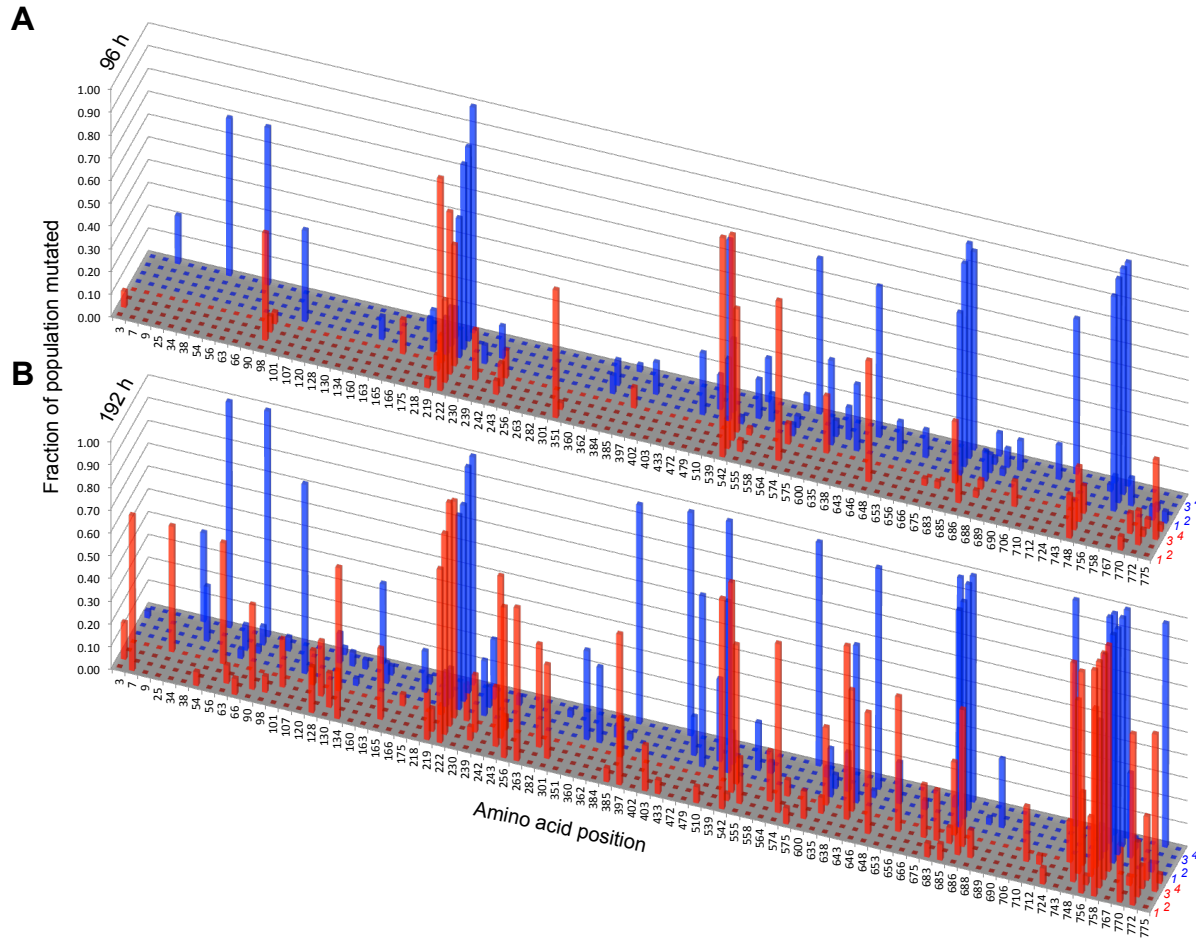


Fig. S9. Mutational frequency of evolved RNA polymerase populations. HTS analysis of mutation frequency at 96 h (A) and 192 h (B). Only mutations present at $\geq 10\%$ at any time point are shown. Red bars represent the four T3 pathway populations while blue bars represent the four SP6 pathway populations.

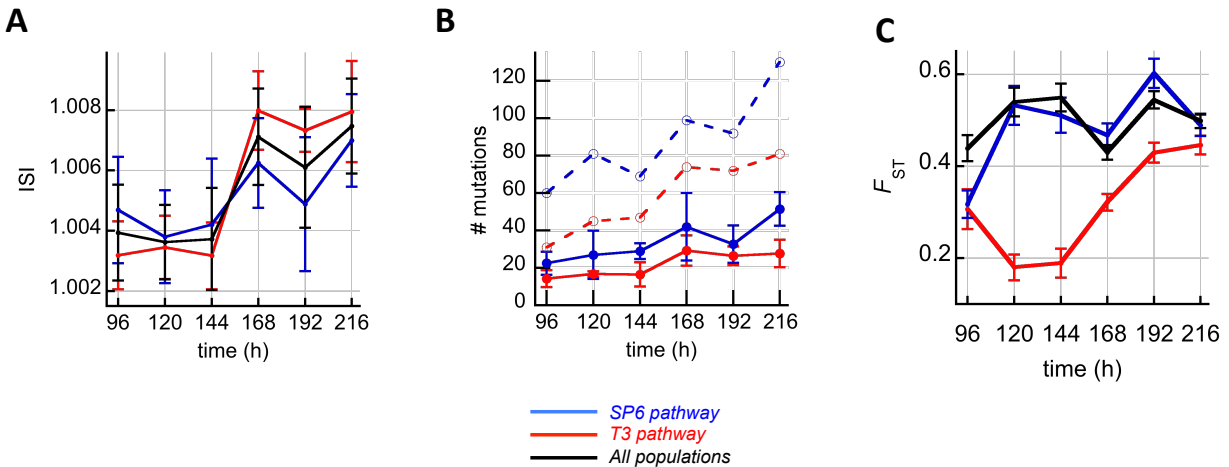


Fig. S10. Time course of genetic diversity and divergence of populations using HTS data. (A) Inverse Simpson index (ISI) diversity, averaged across sequence positions and populations, during the course of the evolution. (B) Average (solid lines) number of unique mutations in each population and total number of unique mutations across all populations (dotted lines) present in $\geq 2.5\%$ abundance during the course of the evolution. Red lines represent the T3 pathway, blue lines represent the SP6 pathway, and black lines represent all populations taken together. Error bars represent the standard deviation among the four populations. (C) F_{ST} values during the course of the evolution. Red populations represent T3 pathway replicates, blue populations represent the SP6 pathway replicates, and black populations represent all evolved populations taken together. Error bars reflect calculated error from high-throughput sequencing (see Supplemental Methods).

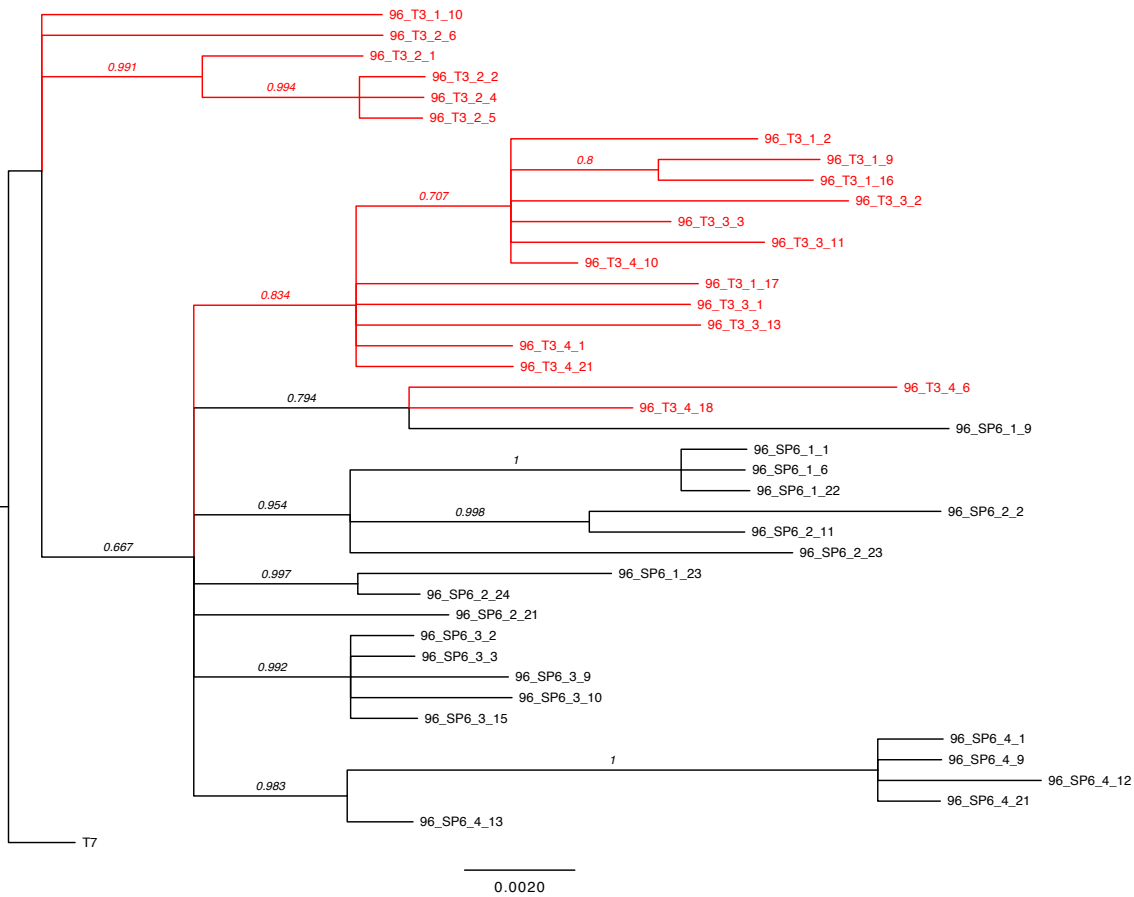


Fig. S11. Phylogenetic analysis of amino acid sequences from several clones from each population at 96 h (see Figs. S4 and S5 for genotypes). Bayesian posterior probabilities are given at the corresponding clade in italics; a polytomy is shown if the posterior probability was < 0.5 . Clones are named with the following convention: **W_XY_Z**, where **W** is the number of hours of evolution (96 or 192 h), **X** is the pathway (T3 or SP6), **Y** is the lagoon replicate, and **Z** is a clone designation. The ancestral sequence (T7) was used as the outgroup. Lineages from the T3 pathway are shown in red; those from the SP6 pathway are shown in black. Scale bar units are substitutions per site.

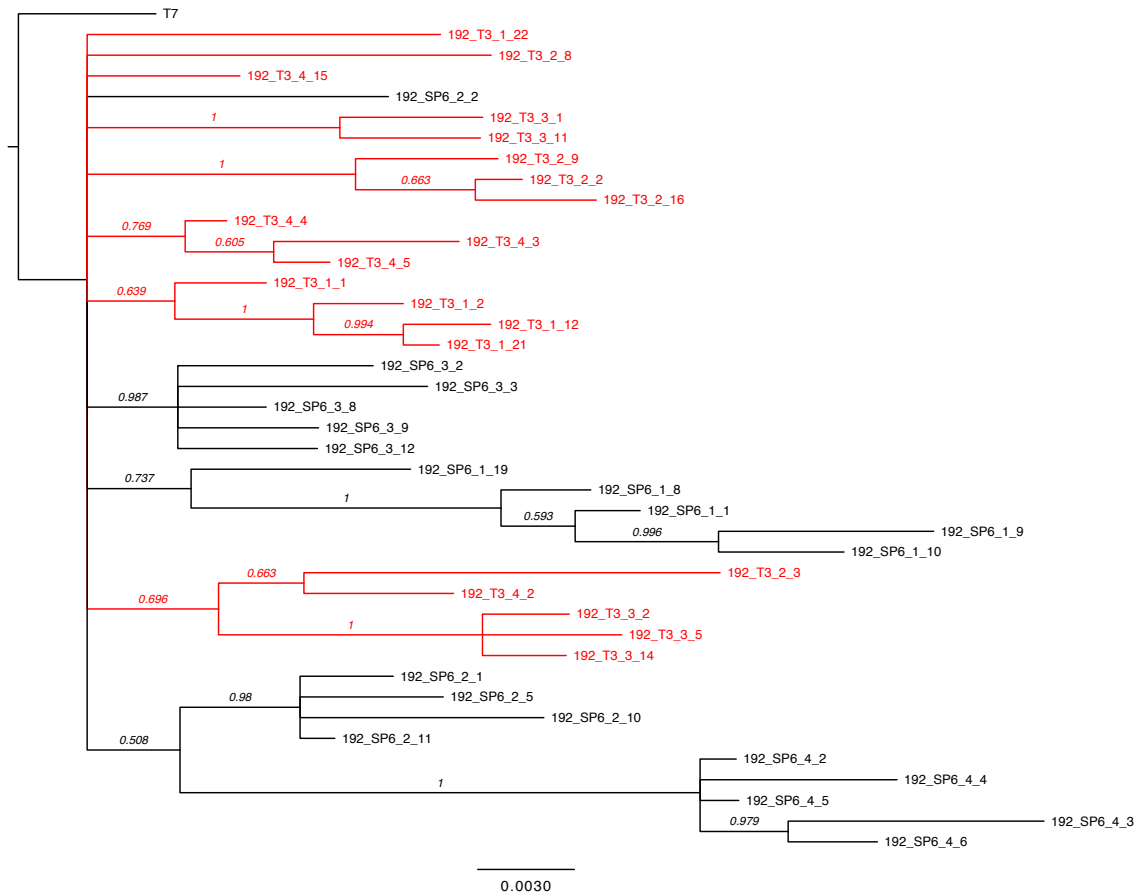


Fig. S12. Phylogenetic analysis of amino acid sequences from several clones from each lagoon at 192 h (genotypes shown in Figs. S6 and S7). Bayesian posterior probabilities are given at the corresponding clade in italics; a polytomy is shown if the posterior probability was < 0.5. Clones are named with the following convention: **W_X_Y_Z**, where **W** is the number of hours of evolution (96 or 192 h), **X** is the pathway (T3 or SP6), **Y** is the lagoon replicate, and **Z** is a clone designation. The ancestral sequence (T7) was used as the outgroup. Lineages from the T3 pathway are shown in red; those from the SP6 pathway are shown in black. Scale bar units are substitutions per site.

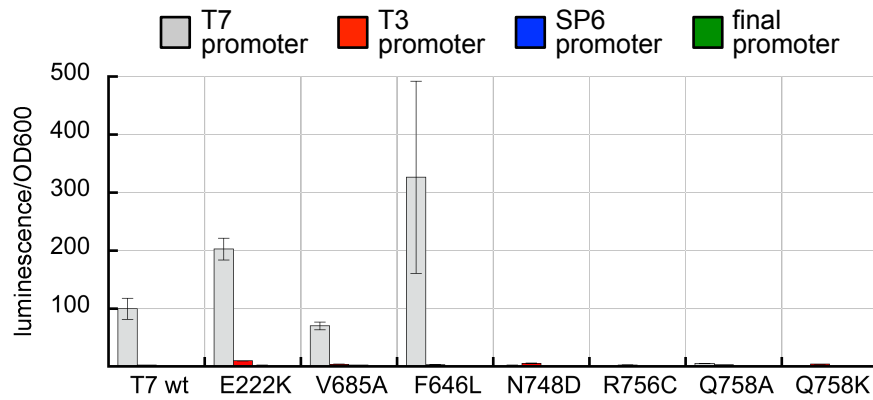


Fig. S13. Activity assays on full panel of promoters of single mutations in T7 RNAP. Full activity profiles for the T7, T3, SP6, and final promoters by in vivo luciferase expression driven by each promoter are shown for wild-type T7 RNAP and wild-type T7 RNAP with each of the single mutations.

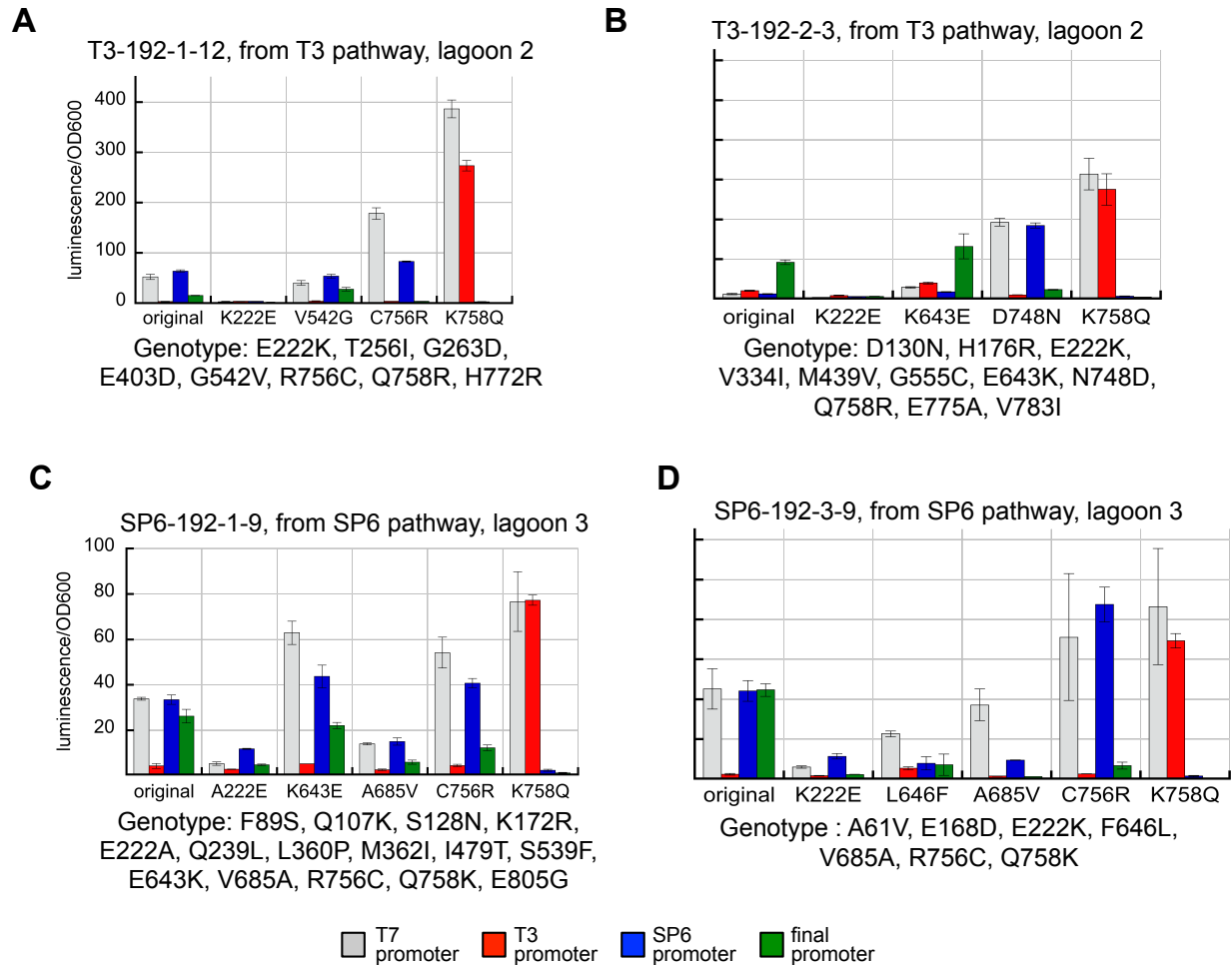


Fig. S14. Activity assays on full panel of promoters of variants from both the T3 pathway and the SP6 pathway with key mutations reverted. Activity profile of two clones from the 192 h time point from the T3 pathway, T3-192-1-12 (A) and T3-192-2-3 (B), and two clones from the 192 h time point from the SP6 pathway, SP6-192-1-9 (C) and SP6-192-3-9 (D), are shown with their full panel of activity assays. Various mutations in each gene were reverted to their wild-type amino acid for each variant and the resulting revertant was assayed on the full panel of promoters.

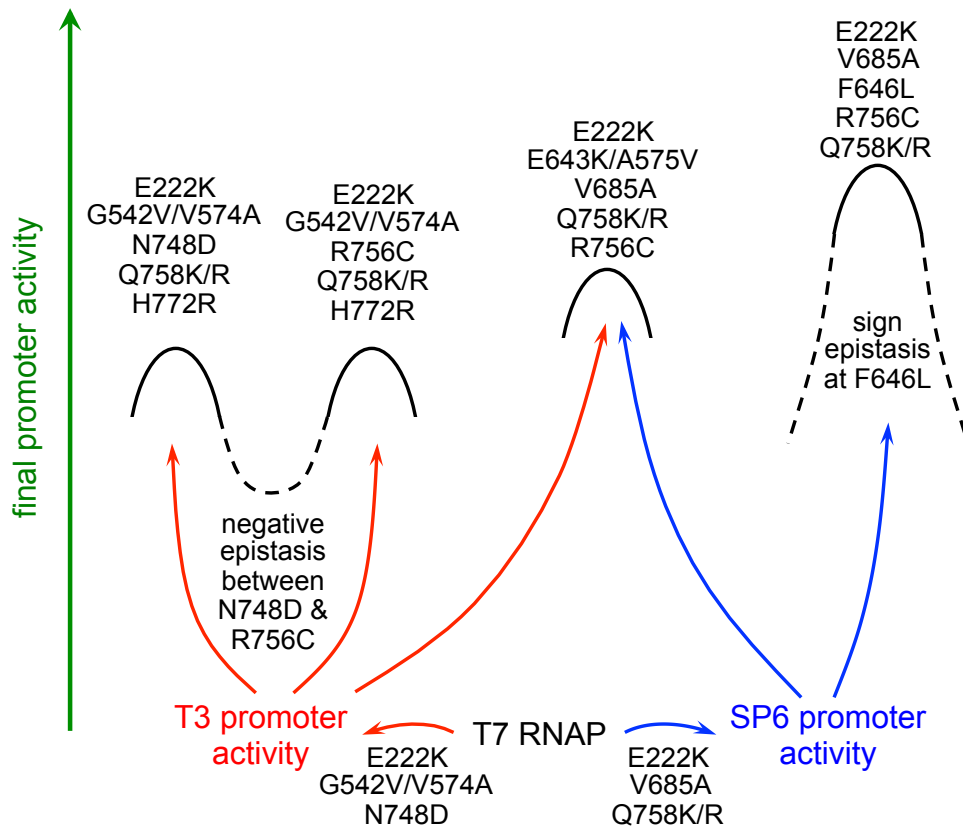


Fig. S15. Diagrammatic representation of examples of fitness peaks and epistasis emerging from both pathways as suggested by genotype and phenotype data. Red arrows represent evolution in T3 pathway populations; blue arrows represent evolution in SP6 pathway populations. Solid peaks represent fitness maxima. Dashed lines represent fitness valleys connecting or surrounding fitness maxima.

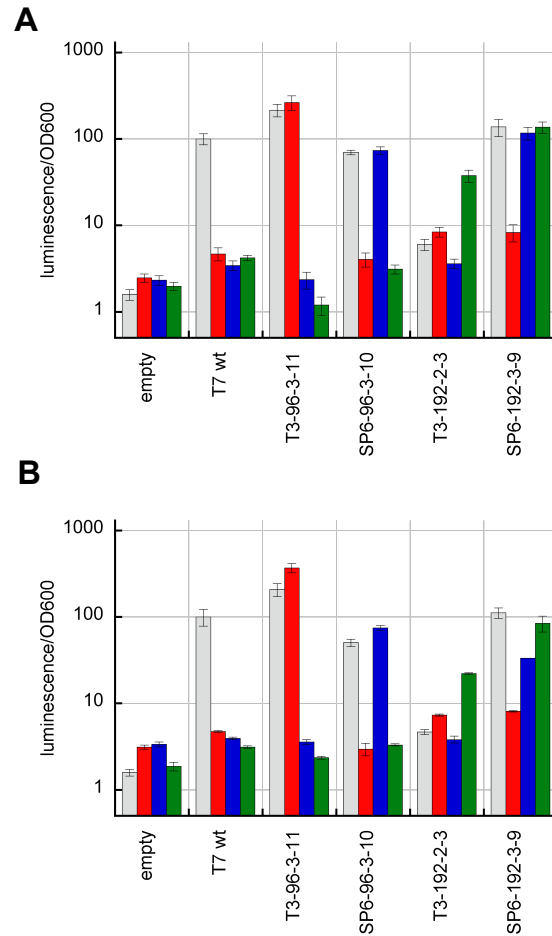


Fig. S16. Comparison of luminescence/OD₆₀₀ measurements of single colonies vs. re-inoculated overnight cultures. (A) Single colonies of transformed cells were picked, grown for 3 h, and then assayed for luciferase activity. (B) Single colonies of transformed cells were picked, grown overnight, then used to re-inoculate 1 mL of media (10 μ L of overnight culture added), grown for 3 h, and then assayed for luciferase activity.