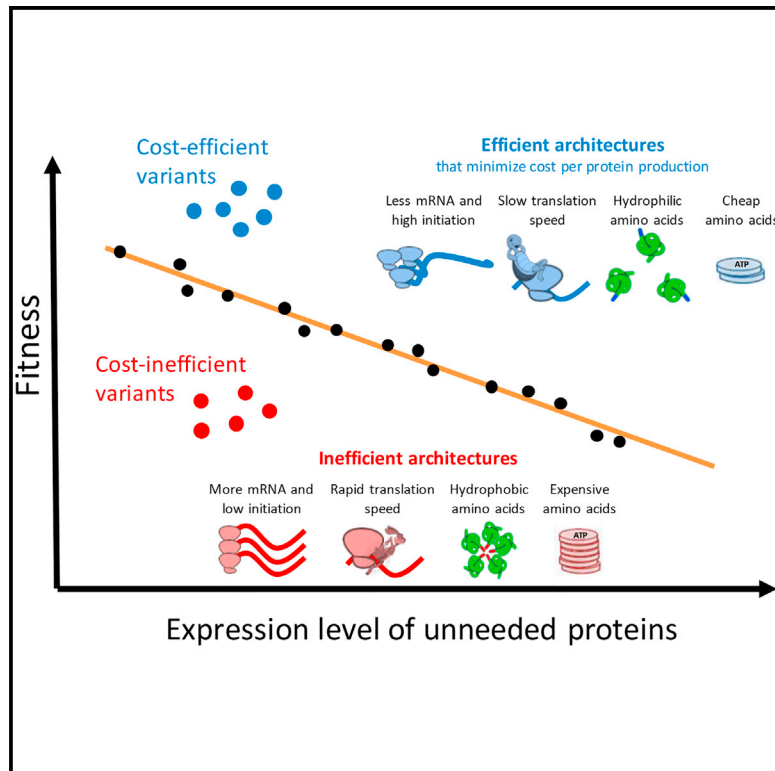# Gene Architectures that Minimize Cost of Gene Expression

## Graphical Abstract



## Authors

Idan Frumkin, Dvir Schirman,
Aviv Rotman, ..., Song Wu,
Sasha F. Levy, Yitzhak Pilpel

## Correspondence

pilpel@weizmann.ac.il

## In Brief

While numerous studies have
investigated regulation of expression
level, Frumkin et al. study gene design
elements that govern expression costs
and allow cells to minimize such costs
while maintaining a given protein
expression level.

## Highlights

- Microorganisms can minimize expression cost with diverse molecular means

- Some design elements can produce more unneeded proteins but maintain high fitness

- Such elements optimize use of production machineries and utilize cheap materials

- Natural highly expressed genes evolved more forcefully to lower expression costs

CellPress

Molecular Cell

# Article

CellPress

# Gene Architectures that Minimize Cost of Gene Expression

Idan Frumkin,[1,5] Dvir Schirman,[1,5] Aviv Rotman,[1,5] Fangfei Li,[2,3] Liron Zahavi,[1] Ernest Mordret,[1] Omer Asraf,[1] Song Wu,[3] Sasha F. Levy,[2,4] and Yitzhak Pilpel[1,6,*]
[1]Department of Molecular Genetics, Weizmann Institute of Science, 7610001 Rehovot, Israel
[2]Laufer Center for Physical and Quantitative Biology
[3]Department of Applied Mathematics and Statistics
[4]Department of Biochemistry and Cell Biology
Stony Brook University, Stony Brook, NY 11794, USA
[5]Co-first author
[6]Lead Contact
*Correspondence: pilpel@weizmann.ac.il
http://dx.doi.org/10.1016/j.molcel.2016.11.007

## SUMMARY

**Gene expression burdens cells by consuming resources and energy. While numerous studies have investigated regulation of expression level, little is known about gene design elements that govern expression costs. Here, we ask how cells minimize production costs while maintaining a given protein expression level and whether there are gene architectures that optimize this process. We measured fitness of ~14,000 *E. coli* strains, each expressing a reporter gene with a unique 5′ architecture. By comparing cost-effective and ineffective architectures, we found that cost per protein molecule could be minimized by lowering transcription levels, regulating translation speeds, and utilizing amino acids that are cheap to synthesize and that are less hydrophobic. We then examined natural *E. coli* genes and found that highly expressed genes have evolved more forcefully to minimize costs associated with their expression. Our study thus elucidates gene design elements that improve the economy of protein expression in natural and heterologous systems.**
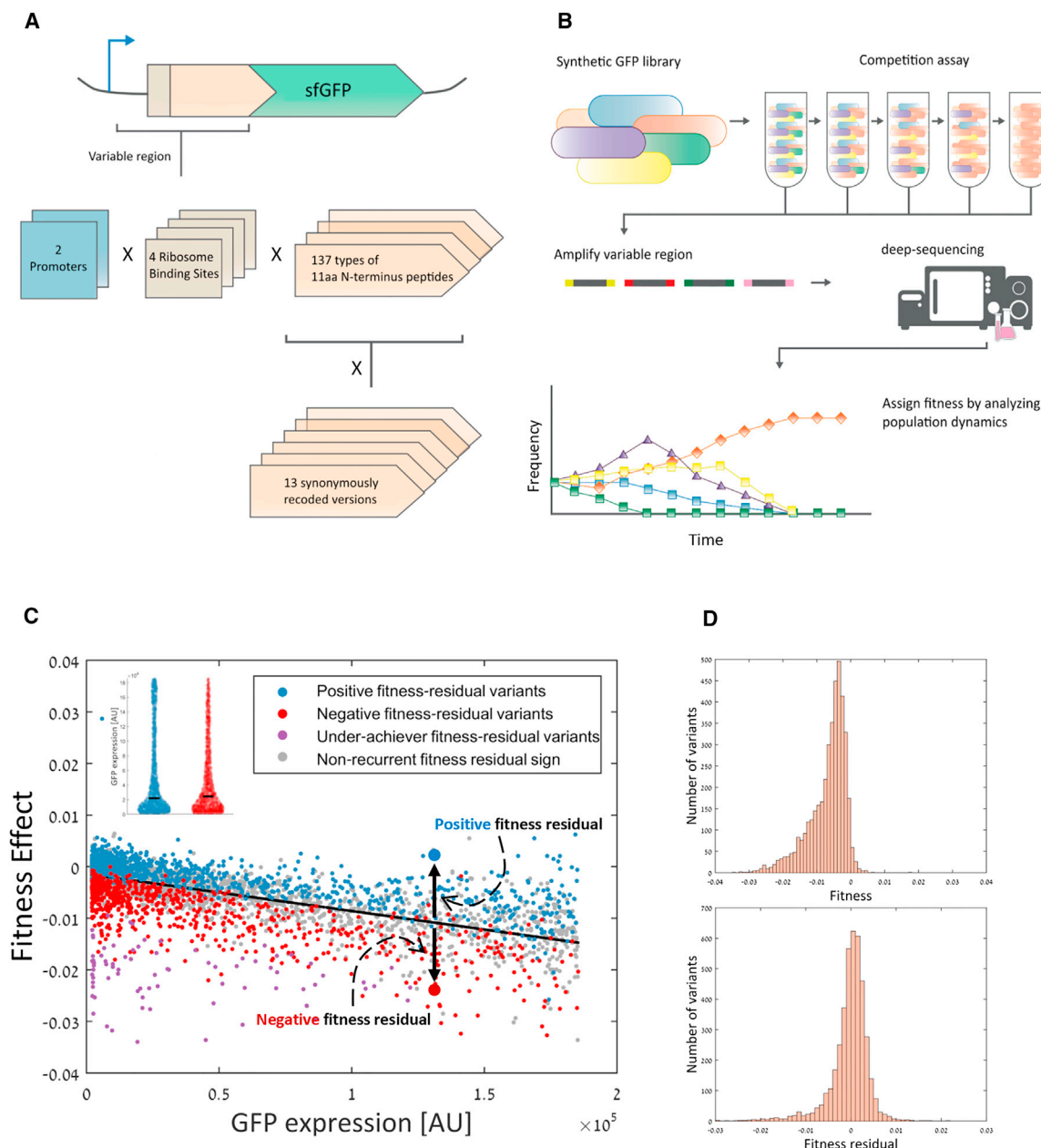
## INTRODUCTION

In nature, cells must express different genes in a regulated manner. On one hand, genes must be expressed at levels that maximize their benefit, and on the other, cells need to minimize the genes' production costs (Dekel and Alon, 2005; Wagner, 2005). Costs of expression originate from spending cellular resources, such as building blocks (amino acids and nucleotides), from allocation of cellular machineries (RNA polymerase and ribosome), and from energy and reducing power consumption (Bienick et al., 2014; Glick, 1995; Ibarra et al., 2002; Rang et al., 2003). Even after their production, proteins might still impose costs when degraded or by exerting toxicity, e.g., due

to aggregation (Geiler-Samerotte et al., 2011). Understanding what molecular processes determine expression cost, its relation to cellular growth and gene regulation, and how costs evolutionarily shape the genome are key aspects of cell biology that remain largely elusive. While numerous studies investigated molecular mechanisms and gene sequence architectures that regulate expression level (Gingold and Pilpel, 2011; Kudla et al., 2009; Qian et al., 2012; Sharp et al., 1986; Subramaniam et al., 2013), very little is known about design elements that govern expression costs.

Different works have studied expression costs in unicellular organisms by imposing the expression of an unneeded protein (Bentley et al., 1990; Dekel and Alon, 2005; Dong et al., 1995; Kafri et al., 2016; Rang et al., 2003; Scott et al., 2010). The production of such unneeded proteins diverts resources from synthesis of the cell's own proteins, thus decreasing cellular fitness (Emilsson and Kurland, 1990; Marr, 1991; Vind et al., 1993). Central to these studies is the characterization of the correlation between the imposed expression levels of the unneeded proteins to the cost. Yet, ultimately natural selection dictates the expression level of natural genes according to the required concentration of each protein. Thus, a fundamental question, which has not been addressed before, is how cells can achieve a specific expression level of a gene while minimizing its expression costs.

Addressing this question is challenging because changes in sequence could affect both expression level and expression costs. To disentangle expression level and expression costs and reveal mechanisms that affect cost per protein molecule, we utilized a synthetic reporter library of ~14,000 different sequence variants, each fused upstream to a GFP gene (Goodman et al., 2013). We then combined competition assays and deep sequencing to measure the fitness of all variants in parallel. This procedure allowed us to elucidate gene architectures that minimize expression cost at a given protein expression level. We show that various molecular mechanisms, such as protein/mRNA ratios, ribosome early elongation pauses, amino acid synthesis costs, and peptide hydrophobicity, determine the cost per protein molecule. We then generated a model that predicts the cost effectiveness of gene architectures and applied it to natural *E. coli* genes. We found that highly expressed genes have

**CellPress**



**Figure 1. 5′ Gene Architectures Affect Cost of Gene Expression at a Given Expression Level**

(A) We utilized a synthetic library of ∼14,000 *E. coli* strains, each expressing a GFP construct with a unique 5′ architecture that includes a promoter, ribosome binding site (RBS), and an 11-amino-acid-fused peptide. There were two different promoter types, four RBSs, and 137 amino acid fusions that were each synonymously re-coded to 13 different versions (see Goodman et al., 2013 for full details).

(B) FitSeq methodology to measure relative fitness of strains in a pooled synthetic library. First, the library was grown six independent times for ∼84 generations, and samples were taken at generations 0, ∼28, ∼56, and ∼84. Then, unique 5′ gene architectures were simultaneously amplified and sent for deep sequencing, which allowed to follow the frequency of each variant in the population over the course of the experiment. Finally, a relative fitness score was assigned for each variant based on its frequency dynamics.

(C) GFP expression level (as measured by Goodman et al., 2013; x axis) versus fitness effect (based on results of repetition C; y axis) of each variant in the library (Pearson correlation, r = −0.79, p < 10$^{-200}$). Fitness effect comes from the burden of expressing unneeded proteins on cellular growth and is calculated by analyzing the frequency dynamics of each variant (see Experimental Procedures). We defined fitness residual as the difference between a variant's observed and expected fitness. The expected fitness is calculated from the regression line between GFP expression and fitness (black line). Some variants consistently demonstrated positive (blue dots, n = 975) or negative (red dots, n = 815) fitness residual sign. Other variants showed extremely low fitness residual, and we termed those variants as "underachievers" (purple dots, n = 80). The group size of positive, negative, and underachiever variants are significantly much higher than expected by chance (Supplemental Information). These results suggest that certain 5′ gene architectures can increase or reduce the cost of gene

*(legend continued on next page)*

**Cell**Press

evolved more forcefully to be encoded by cost-minimizing mechanisms. Our observations indicate that natural selection has shaped genes' architectures to reduce cost of gene expression.

## RESULTS

### 5′ Gene Architecture Affects Cost of Gene Expression

Our question is whether different gene sequence elements can minimize cost of expression per protein molecule and hence increase cellular fitness. To focus on sequence features at the 5′ region of a gene, we utilized a previously published synthetic gene library (Goodman et al., 2013) composed from ~14,000 different variants expressing a GFP gene. Each variant holds a unique variable 5′ gene architecture that includes a promoter, a ribosome binding site (RBS), and an 11-amino-acid-long N terminus fusion (Figure 1A; Experimental Procedures).

To reveal the expression cost of each variant, we measured relative fitness of all variants in parallel in a competition assay in six independent repeats. We then deep sequenced the variable region of the pool of variants and calculated relative fitness of each variant (Figure 1B; see Experimental Procedures).

We regressed fitness values against GFP expression levels and observed a negative, linear correlation (Figure 1C, Pearson correlation, r = −0.79, p < 10$^{-200}$; Figure S1A). The linear decline in fitness with expression is in agreement with previous studies (Kafri et al., 2016; Scott et al., 2010). The regression line, which outlines the relations between fitness and expression, allowed us to estimate the expected fitness for each library variant according to its GFP expression level. Variants whose fitness does not deviate consistently across repeats from this regression line are deduced not to utilize mechanisms that enhance or reduce the production cost per protein molecule.

Yet, many variants did deviate from the linear regression line, demonstrating fitness that is higher or lower than expected given their GFP expression levels. We hypothesized that variants that repeatedly deviated from the expected fitness might utilize gene architectures that either reduce or increase the cost of GFP production per protein molecule. Hence, we calculated each variant's "fitness residual," which we defined as the difference between the actual fitness that we measured for the variant and the fitness expected for it according to its GFP expression level and the linear regression (Figure 1C). A positive fitness residual means that a given variant showed higher fitness than expected given its GFP expression level, suggesting that it can produce this GFP level with lower costs. A negative fitness residual means that the variant showed lower fitness than expected given its GFP expression level.

We then classified each variant as either positive or negative according to its fitness residual sign (Figure 1C, blue and red dots; see Experimental Procedures). Since the observed fitness residual is sensitive to biological noise (i.e., drift during competi-

tion) and experimental errors (i.e., sampling errors), we only classified variants as positive or negative if their fitness residual sign was identical in at least five out of the six repeats of the experiments in each of the two final sampling points of the competition (see Experimental Procedures and Supplemental Experimental Procedures). This approach resulted in 975 positive and 815 negative variants (significantly higher than expected by chance even at very high levels of measurement errors; Supplemental Experimental Procedures). Classification into either positive or negative fitness residual groups allowed us to eliminate the effect of GFP expression level on fitness as these two groups demonstrate the same expression distribution (Figure 1C, inset).

We also noticed a set of 80 library variants, which we termed "underachievers," whose fitness residual scores were repeatedly at the bottom 5% of the entire library (Figure 1C, purple dots; see Experimental Procedures). We hypothesized that these underachiever variants show extremely low fitness residuals because they produce GFP even more wastefully, and we expected them to show stronger usage of low-efficiency gene architectures compared to the negative fitness residual group. There appeared to be no "overachievers" in these data.

### Production of More Proteins per mRNA Molecule Is an Economic Means to Minimize Expression Costs
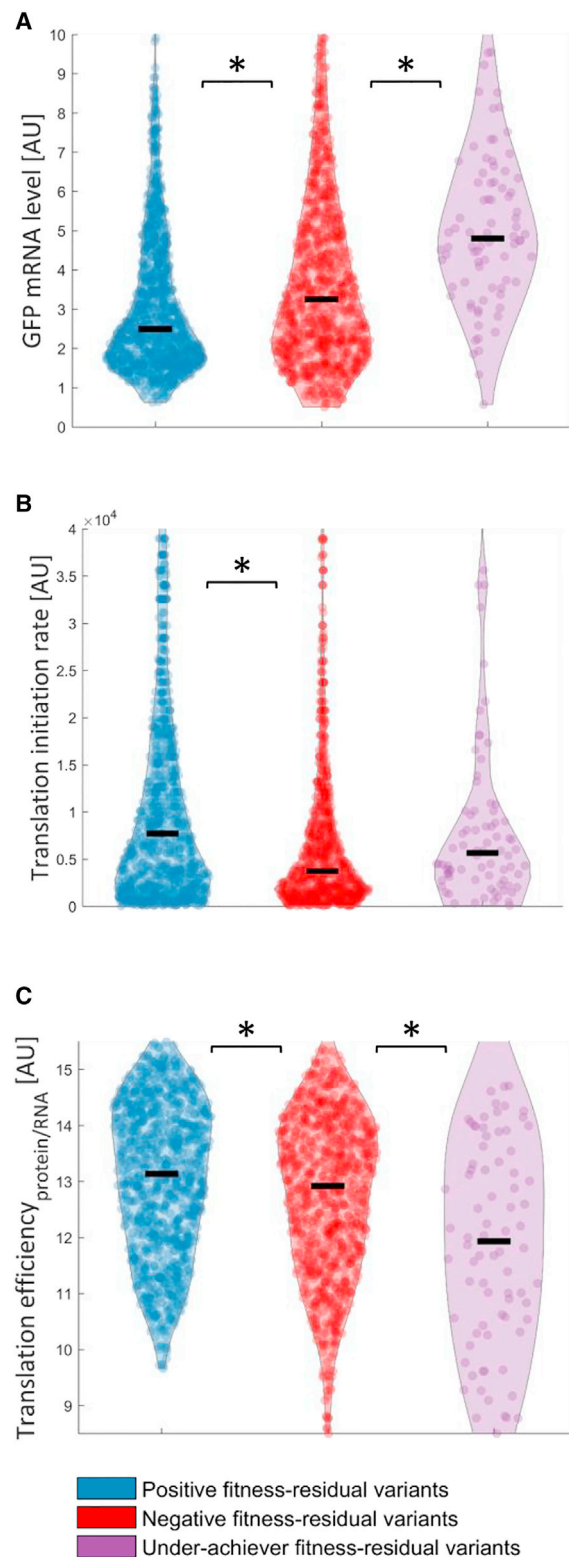
We first hypothesized that reaching the same GFP level with lower levels of mRNA of the GFP gene could be beneficial. While positive and negative fitness residual variants come from the same distribution of GFP expression levels (Figure 1C, inset), we compared their GFP mRNA levels and found positive variants to have lower levels compared to negative variants (Figure 2A; Wilcoxon rank-sum, p = 1.6 × 10$^{-9}$, effect size = 58.26%; see Experimental Procedures). This difference was independent of GFP level: binning the data according to GFP levels, we observed the reduced levels of mRNA for positive variants in all expression bins (Figure S1B).

The observation that positive variants have equal GFP protein levels but lower GFP mRNA levels indicates that they are able to produce more GFP proteins per mRNA molecule. We postulated that high translation initiation rate could be a mechanism for maintaining the same GFP levels despite low mRNA levels in positive variants. We calculated initiation rates for all library variants using the "Ribosome Binding Site Calculator" (Salis, 2011) and observed that indeed positive variants had higher initiation rates (Figure 2B; effect size = 61.9%, Wilcoxon rank-sum, p = 3.7 × 10$^{-18}$). This observation holds true when examining mRNA level versus translation initiation rate at the individual variant level (Figure S2A). Indeed, when examining translation efficiency per variant (using measured protein levels divided by mRNA levels), positive variants demonstrated higher translation efficiencies than negative fitness residual variants (Figure 2C; effect size = 55.67%, Wilcoxon rank-sum, p = 3.4 × 10$^{-5}$). Moreover, we found that underachiever variants demonstrated even

expression. See also Figure S1A. Inset: positive (blue violin plot) and negative (red violin plot) fitness residual variants come from the same distribution of GFP expression level (Wilcoxon rank-sum, p = 0.46). Black line represents the median value. Thus, the effect of GFP levels on fitness was successfully factored out, thus allowing us to elucidate other molecular mechanisms that tune expression cost at given expression levels.

(D) Fitness and fitness residuals demonstrate different distributions. While most variants showed negative fitness values, fitness residual is more similar to a normal distribution, though with a negative tail.

**CellPress**



Positive fitness-residual variants
Negative fitness-residual variants
Under-achiever fitness-residual variants

**Figure 2. Higher Ratio of GFP Protein/mRNA Minimizes Cost of Gene Expression**

(A) Although coming from the same distribution of GFP levels, positive variants (blue violin plot) demonstrate lower mRNA levels of the GFP gene compared to

higher mRNA levels and lower translation efficiencies compared to the negative variants (Figures 2A and 2C; effect size = 68.04% and 63.06%, Wilcoxon rank-sum, p = 9.6 × 10$^{-8}$ and 1.1 × 10$^{-4}$, respectively). Thus, by increasing translation efficiency, cells reduce transcription costs and hence also cost per protein.
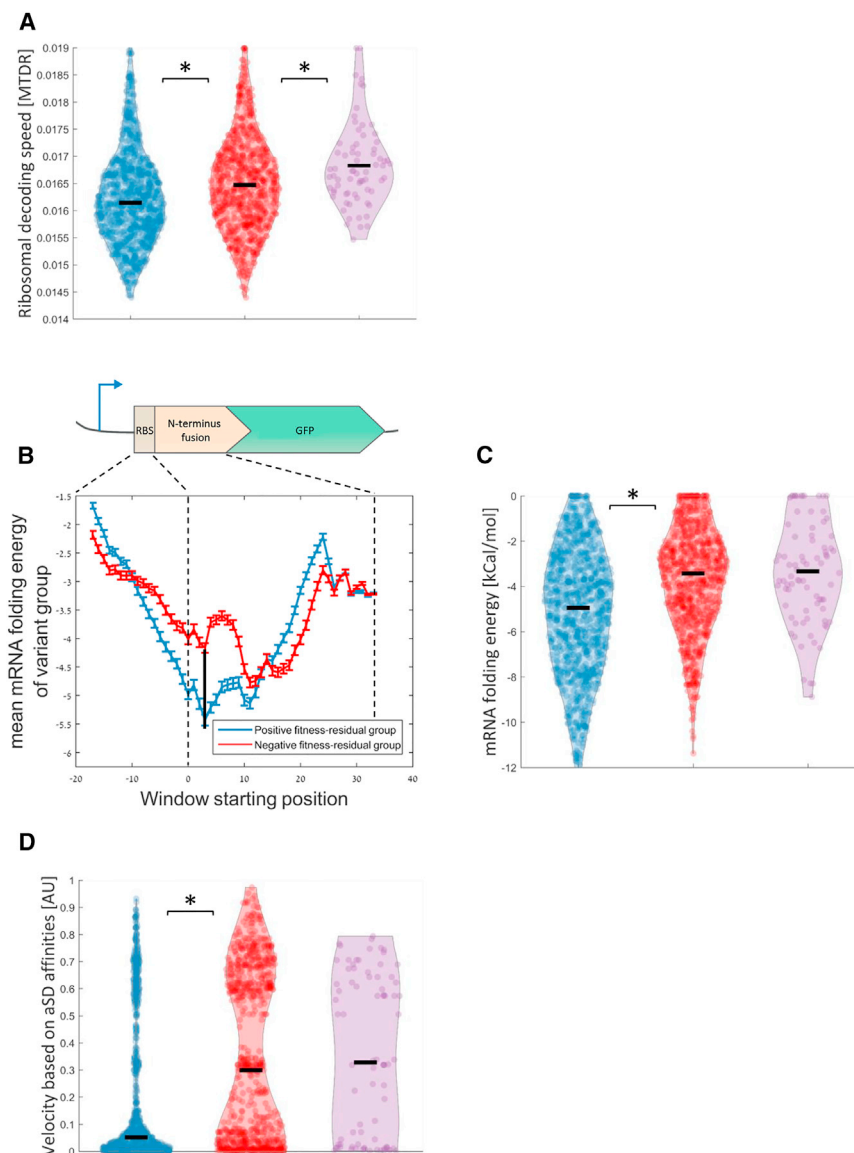
**Slower Translation Speed at Early Elongation of Coding Region, Achieved by Diverse Means, Reduces Expression Costs**

We next aimed to elucidate other cellular mechanisms that directly regulate the translation machinery and that might reduce expression costs. We first examined codon decoding speeds by the ribosome. Codon adaptation of transcripts to the cellular tRNA pool has been shown to be a regulatory mechanism for translation elongation (Goodarzi et al., 2016; Higgs and Ran, 2008; Kudla et al., 2009; Plotkin and Kudla, 2011; Shah and Gilchrist, 2011; Weinberg et al., 2016; Yona et al., 2013). Specifically, the prevalence of slowly translated codons at the 5′ of open reading frames (ORFs) has been suggested to support the efficiency of gene translation (Tuller et al., 2010a). This "ramp model" proposes that delaying ribosomes at the beginning of the elongation phase decreases downstream ribosomal pauses and collisions, which can therefore reduce ribosome jamming, and perhaps also ribosomal abortion events.

Although contradicting evidence were reported for the existence and relevance of this mechanism to expression level (Charneski and Hurst, 2014; Dana and Tuller, 2014; Heyer and Moore, 2016; Ingolia et al., 2009; Shah et al., 2013; Tuller and Zur, 2015), the main prediction of the model—that 5′ ramping reduces cost of expression at a given expression level—has not been tested so far. Here, we had the first opportunity to test this hypothesis as only the 5′ variable region of the GFP varied in the library, while all other parameters remained constant. Thus, we asked whether slow 5′ translation speed is associated with positive fitness residual. We used "mean of the typical decoding rates" (MTDR) (Dana and Tuller, 2014), a measure of codon decoding time derived empirically from ribosome profiling data in *E. coli* (see Experimental Procedures), to calculate translation speed for each library variant. We reasoned that if translational ramp is beneficial, then low MTDR scores, i.e., low ribosome speeds, should be more prevalent among the positive fitness residual variants. Indeed, our results showed that positive variants demonstrate significantly lower translation speeds at the N-terminal fusion (Figure 3A; effect size = 59.55%, Wilcoxon rank-sum, p = 3 × 10$^{-12}$) and further for

negative variants (red violin plot) (effect size = 58.26%, Wilcoxon rank-sum, p = 1.6 × 10$^{-9}$). Consistently, underachiever variants (purple violin plot) show higher mRNA levels compared to negative variants (effect size = 68.04%, Wilcoxon rank-sum, p = 9.6 × 10$^{-8}$). Black line represents the median value. (B) Positive variants show higher translation initiation rates compared to negative variants (effect size = 61.9%, Wilcoxon rank-sum, p = 3.7 × 10$^{-18}$). (C) Positive variants demonstrate higher translation efficiencies (protein/mRNA) compared to negative variants (effect size = 55.67%, Wilcoxon rank-sum, p = 3.4 × 10$^{-5}$). Consistently, underachiever variants (purple violin plot) further show lower translation efficiencies compared to negative variants (effect size = 63.06%, Wilcoxon rank-sum, p = 1.1 × 10$^{-4}$).
Statistically significant differences (p < 0.05) are marked with an asterisk. See also Figures S1B and S2A.

CellPress

**A**



**B**



**C**



**D**



**Figure 3. Slow Translation Speed at Early Elongation, Achieved by Diverse Molecular Means, Reduces Expression Cost**

(A, C, and D) Positive variants show lower values of codon decoding speed (A), stronger mRNA structures (C), and lower speeds due to higher anti-Shine Dalgarno affinities (D) compared to negative variants (effect size = 59.55%, 65.03%, and 63.82%, Wilcoxon rank-sum, p = 3 × 10$^{-12}$, 5.4 × 10$^{-28}$, and 6.3 × 10$^{-24}$, respectively). Statistically significant differences (p < 0.05) are marked with an asterisk. See also Figure S1B.

(B) Mean folding energy of mRNA secondary structure according to window's start position for positive (blue curve) and negative (red curve) variants; error bars represent SEM. Dashed lines mark different positions along the variable region upstream to the GFP. Black vertical line marks the beginning of window with the largest observed difference, which is found at nucleotide positions +4 of the ORF, just after the first AUG codon. The distributions at this window position are seen in (C). See also Figure S2B.

to negative variants along many different window positions (Figure 3B; Figure S2B for different window sizes). Strikingly, the maximum difference in folding energy is observed when the window's start position is at the beginning of the translated region of the ORF, excluding the upstream 5′ UTR (Figure 3C; effect size = 65.03%, Wilcoxon rank-sum, p = 5.4 × 10$^{-28}$). Hence, these results, together with previous ones, reveal the dual role of mRNA folding: on one hand, loose mRNA structure at the RBS is associated with high expression level (Goodman et al., 2013), and on the other hand, utilization of a strong secondary structure at the 5′ end of the ORF can reduce per-protein costs.
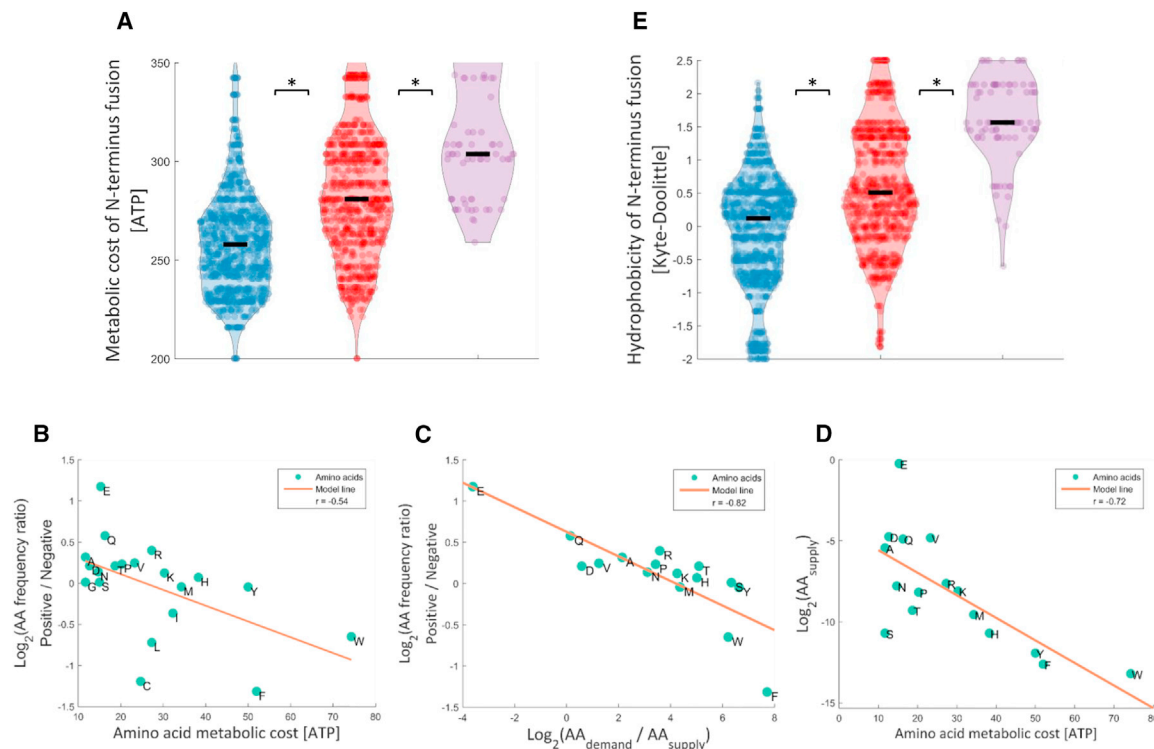
the underachievers (effect size = 64.79%, Wilcoxon rank-sum, p = 1.2 × 10$^{-5}$).

Though in the original ramp model ribosome attenuation was proposed to be obtained by codons that correspond to rare tRNAs, additional mechanisms that can slow down the ribosome at early elongation regions could serve in ramping. These mechanisms include, in particular, tight mRNA secondary structure (Goodman et al., 2013; Tholstrup et al., 2012; Tuller et al., 2010b; Wen et al., 2008) and high affinity to the anti-Shine Dalgarno (aSD) motif of the ribosome (Li et al., 2012). We thus examined each of these factors separately and asked whether they are associated with positive or negative fitness residual.

When we computed folding energies for segments of mRNA nucleotides on a sliding window along the variable region of each variant, we found that positive fitness residual variants demonstrated tighter secondary structures compared

It was previously suggested that elongating ribosomes in *E. coli* dwell longer on sequences that have high affinity to the aSD motif in the ribosome (Li et al., 2012). However, this observation has been recently questioned (Mohammad et al., 2016). We next examined the effects of Shine Dalgarno-mediated ribosomal pauses on fitness residuals. We calculated affinities to the aSD along the sequence of each variant, derived a ribosome speed estimation based on these affinities (see Experimental Procedures) and found that positive fitness residual variants are characterized by low ribosome speed early in the ORF (Figure 3D; effect size = 63.82%, Wilcoxon rank-sum test, p = 6.3 × 10$^{-24}$).

We thus provide the first experimental evidence for a set of three gene architecture factors—codon decoding time, mRNA structure, and affinity to the anti-Shine Dalgarno motif—that could each implement 5′ ramping by slowing down ribosomes and, by that, allow cells to reduce the cost of gene expression at a given expression level.

**Figure 4. Usage of Expensive-to-Synthetize, Lowly Available, and Hydrophobic Amino Acids Decreases Fitness Residual**

(A) N terminus amino acid fusions of negative variants are more expensive to synthesize compared to positive variants (effect size = 72.74%, Wilcoxon rank-sum, p = 7.4 × 10⁻⁶². Underachievers utilize even more expensive amino acids (effect size = 72.75%, Wilcoxon rank-sum, p = 1.7 × 10⁻¹¹). See also Figures S1B and S2C.

(B) The frequency ratio of amino acids between positive and negative variants is negatively correlated with the energetic cost of amino acids (Pearson correlation, r = −0.54, p = 0.01). Each amino acid is marked according to its one-letter code.

(C) The frequency ratio of amino acids between positive and negative variants is negatively correlated with the demand/supply ratio of amino acids (Pearson correlation, r = −0.82, p = 10⁻⁴). Demand comes from occupancy of ribosomes on each transcript (see Experimental Procedures), and supply is the cellular concentration of each amino acid (Bennett et al., 2009).

(D) Amino acid availability and energetic cost are correlated (Pearson correlation, r = −0.72, p = 1.8 × 10⁻³).

(E) N terminus amino acid fusions of negative variants are more hydrophobic than positive variants (effect size = 69.11%, Wilcoxon rank-sum, p = 3.2 × 10⁻⁴⁴). N terminus fusion of underachievers are even more hydrophobic (effect size = 81.67%, Wilcoxon rank-sum, p = 7.7 × 10⁻²¹). See also Figures S1B and S2C.

Another means of reducing translation speed that was recently demonstrated (so far in yeast) is the incorporation of positively charged amino acids (Charneski and Hurst, 2013) or proline residues (Artieri and Fraser, 2014) in newly synthesized peptides. Yet, we did not detect any difference in frequency of such amino acids between the positive and negative fitness residual groups.

### Amino Acid Synthesis Cost and Hydrophobicity Affect Cost of Gene Expression

So far we have examined features that are based on the nucleotide sequence and how it associates with fitness residual. Next, we aimed to explore the possibility that the amino acid composition of the N terminus fusion to the GFP associates with cellular fitness.

Amino acids differ by the metabolic costs associated with their biosynthesis—predominantly energy and reducing power determinants invested in their metabolic production (Akashi and Gojobori, 2002). We thus hypothesized that usage of energetically expensive amino acids may cause a heavier burden at a given

expression level. Indeed, lower cost of the N terminus fusions were found to associate with positive fitness residual variants (Figure 4A; effect size = 72.74%, Wilcoxon rank-sum, p = 7.4 × 10⁻⁶². Here, as well, underachiever variants show more expensive amino acid usage compared to the negative group (Figure 4A; effect size = 72.75%, Wilcoxon rank-sum, p = 1.7 × 10⁻¹¹).

We further examined the relation between fitness residual and amino acid energetic cost by calculating the frequency ratio of each individual amino acid between the positive and negative fitness residual groups (see Experimental Procedures). Remarkably, this frequency ratio was found to negatively correlate with the metabolic cost of each amino acid (Figure 4B; Pearson correlation, r = −0.54, p = 0.01). These observations suggest that expensive-to-synthesize amino acids burden cells during their costly production due to a potential feedback that increases their synthesis in response to consumption.

In addition to direct metabolic cost, the incorporation of amino acids that appear in low cellular concentrations could reduce

fitness indirectly as it might disturb the synthesis of other native proteins. We used ribosome profiling data (Li et al., 2012) to calculate amino acid demands and utilized previously measured cellular concentrations as amino acid supplies (Bennett et al., 2009) (see Experimental Procedures). Indeed, we found that amino acids with low demand-to-supply ratios are more prevalent in positive variants (Figure 4C; Pearson correlation, $r = -0.82$, $p = 10^{-4}$). This observation implies that utilization of amino acids that are less available to the cell (either due to high demand or low supply) increase expression cost and are associated with negative fitness residual variants. Since metabolic cost of amino acids and their cellular supplies are correlated (Figure 4D; Pearson correlation, $r = -0.72$, $p = 1.8 \times 10^{-3}$), we could not evaluate which mechanism—cost or availability—contributes more to fitness residual.

We next reasoned that an additional factor by which a protein could affect fitness is its toxicity, e.g., due to aggregation. As aggregation is driven by hydrophobic interactions, we turned to a conventional measure of amino acid hydrophobicity (Kyte and Doolittle, 1982) to examine whether it is predictive of fitness residuals. We found that positive fitness residual variants tended to have significantly less hydrophobic amino acids fused to the GFP (Figure 4E; effect size = 69.11%, Wilcoxon rank-sum, $p = 3.2 \times 10^{-44}$). Underachievers showed an even more pronounced effect (Figure 4E; effect size = 81.67%, Wilcoxon rank-sum, $p = 7.7 \times 10^{-21}$). This negative effect of hydrophobic residues in cytosolic proteins could indeed be derived from post-synthesis costs, but it could also reflect an equally interesting possibility: that aggregation-prone peptides reduce the functional level of the GFP (and similarly the fraction of the active form of native proteins). According to this possibility, aggregation is wasteful and must be compensated by further costly production to reach the required expression level of the protein.

We further found that the higher the GFP expression, the more beneficial it should be to utilize cheap or hydrophilic amino acids (Figure S2C).

### All Sequence Parameters Contribute Independently to Fitness

We have revealed, so far, a set of mechanisms that affect expression costs and therefore cellular fitness. Although these mechanisms are different in their nature, it is possible that variants that score highly on one of these parameters tend to score highly on others. For example, anti-Shine Dalgarno affinity could correlate with the energy of the secondary structure of the mRNA, as both parameters are influenced by Guanine content. To check this possibility, we computed the correlation among the variants in the library between each pair of sequence parameters: codon decoding speed, mRNA secondary structure, anti-Shine Dalgarno affinity, hydrophobicity, and amino acid energy cost. Reassuringly, no strong correlation was found between any two parameters (Figure 5). Nonetheless, for feature pairs that did demonstrate non-negligible correlations (Pearson correlation, $r > 0.1$), we asked whether the signal of one feature is still observed while controlling for variation in the other. We found that each factor contributed directly to the signal, even upon

controlling for other factors as potential confounders (see Figure S3).

### Expression Costs Can be Minimized Even at Specified Amino Acid Sequences

Since maintaining a protein's function usually requires keeping its specific amino acid sequence, we next asked whether the mechanisms that we found here can reduce expression costs for a specified peptide sequence by using alternative nucleotide sequences. We defined "Δfitness-residual" as the difference between a variant's fitness residual and the average fitness residual of all library variants who share with that variant the same amino acid sequence. Then, we compared the various architectural features between variants with above-average Δfitness-residual to variants with below-average Δfitness-residual (see Experimental Procedures).
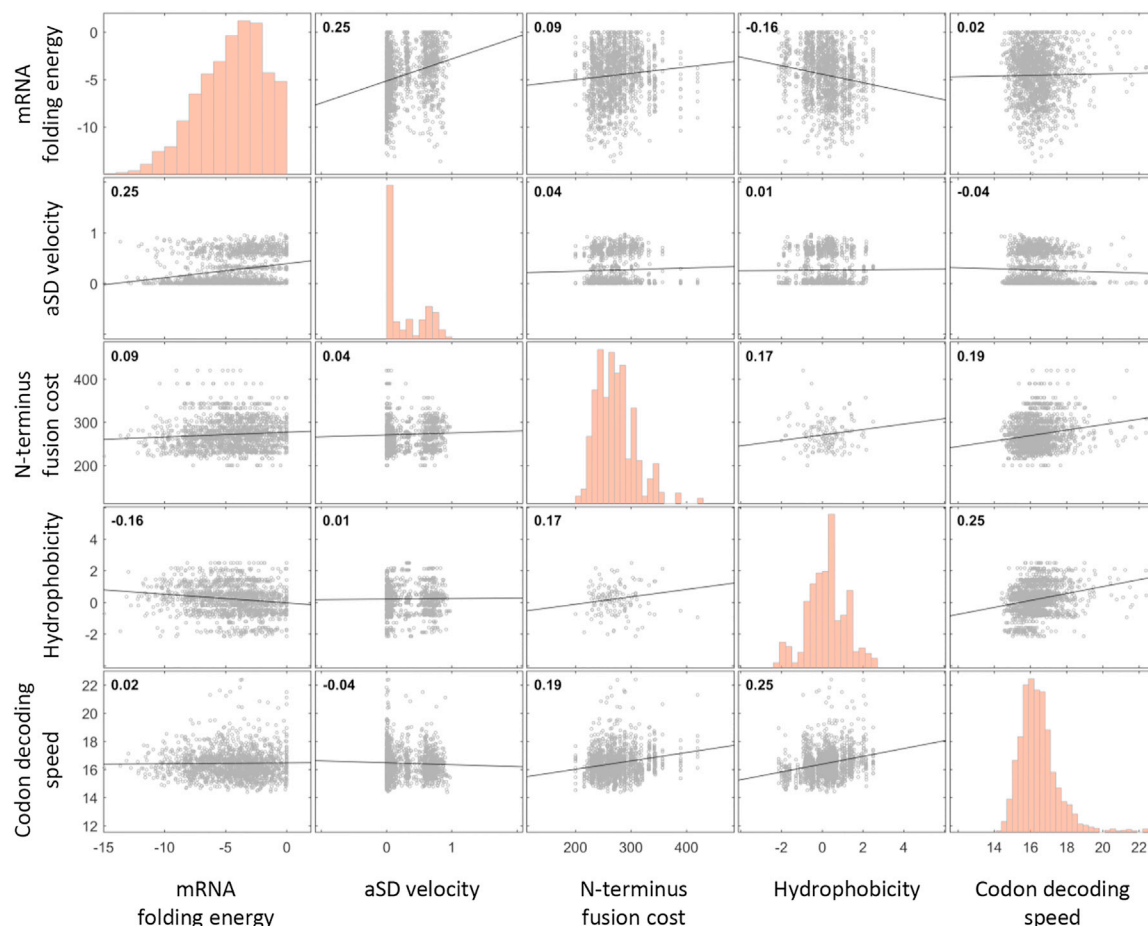
Figures 6A–6E depict, for each of the analyzed features, the difference in feature value between variants with above- or below-average Δfitness-residual. Interestingly, for each feature, the above- and below-average sub-groups had significantly different feature scores, reflecting the same trends as observed in all earlier analyses. For example, mRNA levels tend to be higher in the below-average sub-group in most of the 137 N terminus fusions (t test, p values for GFP mRNA levels = $6.2 \times 10^{-3}$, initiation rates = $7 \times 10^{-9}$, codon decoding speeds = $4.3 \times 10^{-2}$, mRNA folding = $3.5 \times 10^{-16}$, and aSD velocity = $7.6 \times 10^{-7}$). The conclusion from this analysis is that although amino acid features affect fitness residuals, the other features provide sufficient degrees of freedom to minimize costs even at a specified amino acid sequence.

### A Regression Model Calculates Relative Contribution of Each Feature and Predicts Fitness Residual Scores

So far, we have examined fitness residual as a binary classification, namely categorizing variants with either positive or negative fitness residual. Complementing this binary analysis, in Figure S4A, we show that each feature correlates significantly with actual fitness residual values. We next aimed to predict actual fitness residual values of the library variants from their gene architecture features using a multiple linear regression model. We trained the model on a randomly chosen subset of 70% of the library variants, cross validated it on all other variants by comparing their predicted and observed fitness residual, and found a good correlation (see Experimental Procedures; Figure 7A; $r = 0.53$, $p < 10^{-200}$).

When the regression was performed on a scrambled library, which randomly links feature values and variants, the correlation between observed and predicted fitness residual was practically eliminated (Figure S4B; $r = 0.02$). We performed $10^5$ such randomizations, and all of them demonstrated such extremely weak correlations. This negative control demonstrates that we obtained a genuine means to predict fitness residual values based on computable gene architecture parameters. We concluded that a gene architecture that utilizes more of the features that we discovered and that, to a greater extent, typically gives rise to higher fitness residuals as expression costs are further minimized.

**CellPress**



**Figure 5. Each Feature Affects Fitness Residual Independently**

Correlation plots of each feature pair show lack of correlation in most cases and only weak correlations in other cases. For feature pairs with Pearson correlation of r > 0.1, we compared the difference in one feature while controlling for the second and vice versa. See also Figure S3. Black lines are the regression curves between each feature pair. Number at upper-left corner is the Pearson correlation.
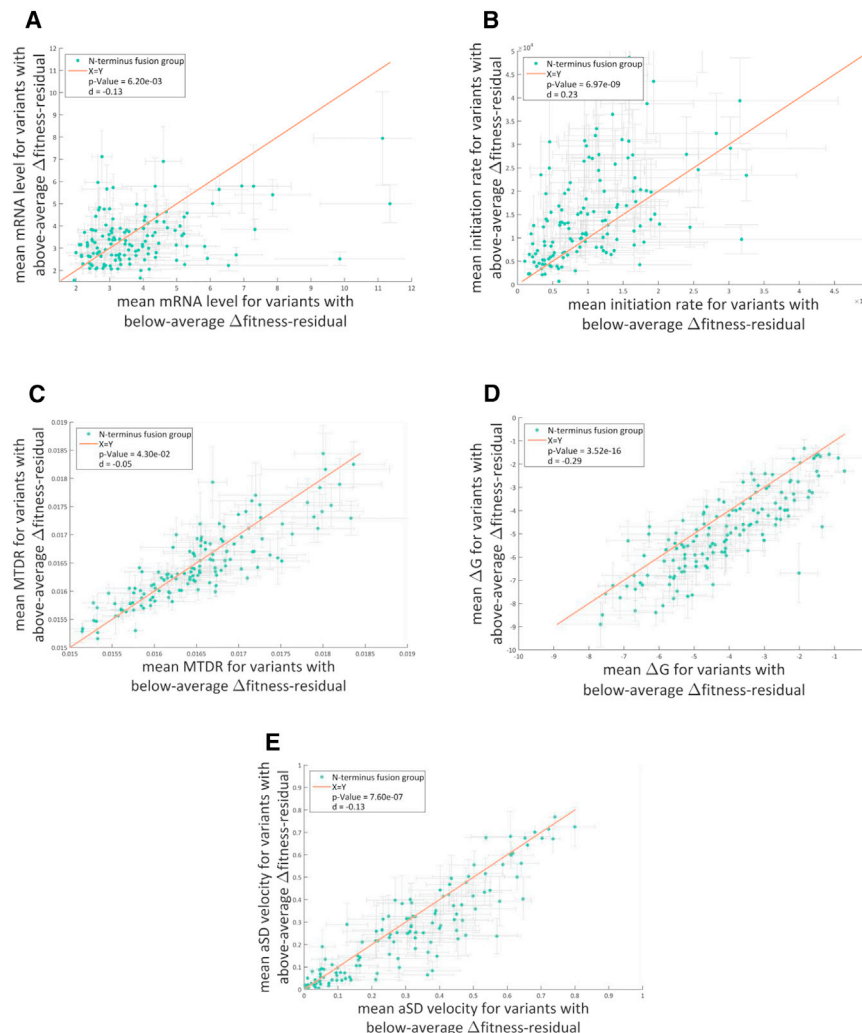
Additionally, this regression model allowed us to calculate the relative contribution of each feature by comparing the coefficients assigned by the regression model (Figure 7B). This analysis revealed that the features contributing to fitness residual the most are hydrophobicity and metabolic cost of the N terminus fusion, while codon decoding speed contributes the least. To avoid over-fitting of our model on the library data, we performed feature selection using the Lasso algorithm (see Experimental Procedures). This validation resulted in the exclusion of only codon decoding speed from the model, suggesting that its contribution to fitness residual is indeed lower compared to other features.

**Highly Expressed Natural Bacterial Genes Have Evolved Gene Architectures that Minimize Their Production Costs**

With these findings from the synthetic library, we next asked whether the mechanisms that we revealed as cost reducing were also utilized by natural selection to optimize *E. coli*'s native genes. We thus calculated each *E. coli* gene's score

with respect to the relevant features and used the regression model to predict its fitness residual score (see Experimental Procedures and Table S4, related to Figure 7). Since a higher expression level results in higher expression cost, we next hypothesized that *E. coli* genes with higher expression levels are more likely to be endowed with cost-reducing architectures. Indeed, we found a significant correlation between predicted fitness residual of *E. coli* genes and their protein expression levels (Figure 7C; r = 0.25, p = $2 \times 10^{-53}$), demonstrating a stronger selection for optimizing the 5′ gene architecture for highly expressed genes. We obtained similar results when predicting fitness residuals for all genes in the Gram positive *B. subtilis*, pointing to the generality of the model (Figure 7E; r = 0.33, p = $10^{-93}$; see Experimental Procedures and Table S4, related to Figure 7).

Interestingly, the range of fitness residuals predicted by our model for the *E. coli* and *B. subtilis* genes was significantly larger than the range predicted by a mock regression model that was trained on randomly scrambled data of the synthetic library (see Experimental Procedures; Figures 7D and 7F; p < $10^{-5}$).

**Figure 6. Variant with Same N Terminus Amino Acid Fusion Demonstrate a Range of Fitness Residuals**

(A–E) Each dot represents one of the 137 N terminus fusions in the library. The x axis and the y axis represent the mean value of a feature for the variants with either below-average or above-average Δfitness-residual, respectively. The vertical and horizontal error bars represent standard errors for each of the axes. A statistical difference for deviance from the X = Y line was observed for all features, suggesting that even at a given amino acid sequence, these mechanisms affect fitness residual and can minimize expression costs (t test, p values: A, mRNA levels, $6.2 \times 10^{-3}$; B, initiation rates, $7 \times 10^{-9}$; C, codon decoding speeds, $4.3 \times 10^{-2}$; D, mRNA folding, $3.5 \times 10^{-16}$; and E, aSD velocity, $7.6 \times 10^{-7}$). d is Cohen's d that calculates the effect size.

This observation suggests that the model that we trained on the library data is able to expose the expression-cost optimality of natural 5′ gene architectures.
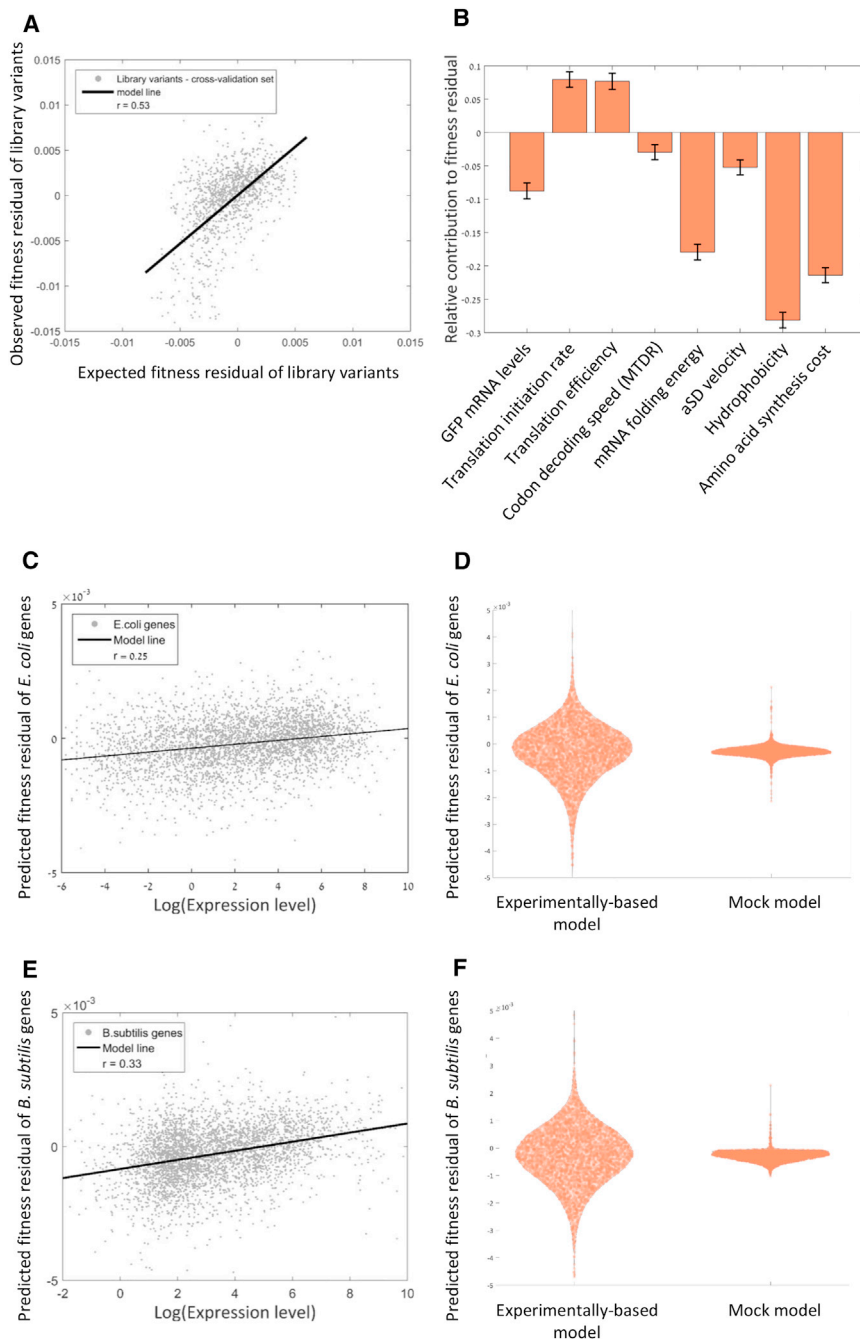
## DISCUSSION

In this study, we found architectures and motifs that govern expression costs and reveal their function even beyond a direct effect on the process of expression. We show that regulating initiation and mRNA levels affects expression cost, as increasing the number of proteins that are produced per mRNA is associated with a positive fitness residual. This architecture could be beneficial because it reduces energy and resource consumption that are devoted to mRNA production. If cost reducing, why do genomes not further utilize the strategy of low transcription and mRNA abundance, combined with high translation initiation? One potential reason is that too low of mRNA levels might lead to increased expression noise (Taniguchi et al., 2010) or increased response time to an environmental signal (Gasch and Werner-Washburne, 2002). It is thus expected that natural genes would show a tradeoff between cost-reducing architec-

tures and designs that satisfy other requirements, such as controlled noise and short response times.

The "translational ramp" theory predicted an effect of ribosome speed at early elongation on expression cost at a given expression level (Tuller et al., 2010a). The theory was never tested as such, since fitness reduction upon expression of an unneeded protein was not systematically measured for different gene sequences at various expression levels. We demonstrate here that slow translation speed at the 5′ end is beneficial in terms of reduced expression cost and increased cellular growth rate. We show that in addition to codon decoding times, there are at least two additional ramping means that are likely beneficial: occurrence of Shine-Dalgarno-like sequences and strong secondary structures.

Recent works showed that 5′ mRNA secondary structure governs expression level of transcripts in bacteria (Goodman et al., 2013; Kudla et al., 2009; Shah et al., 2013). Here, we observed that tight mRNA structures are enriched in positive variants. Consequently, it seems that mRNA structure plays a more complex role than previously thought. On one hand, 5′ mRNA structure, specifically upstream of the AUG start codon, regulates expression levels as it governs initiation rates (Goodman et al., 2013; Salis, 2011). On the other hand, tight structures at the beginning of the ORF, which were previously observed in *E. coli* genes (Tuller et al., 2011), are shown here to be beneficial in minimizing expression cost.

We revealed that the amino acid composition of a gene can also affect expression cost at a given expression level by showing that hydrophobic amino acids reduce fitness residual, perhaps due to their increased tendency to form toxic aggregates in the cytoplasm. In agreement with this, it was shown that mis-folded proteins impose growth reduction to yeast

**CellPress**



**Figure 7. A Model that Predicts Fitness Residual Accurately Reveals that Fitness Residual of Natural Bacterial Genes Is Correlated with Their Expression Level**

(A) A linear regression model based on all eight features predicts fitness residual accurately in a cross-validation test (Pearson correlation, r = 0.53, p < $10^{-200}$). See also Figure S4.

(B) The weighted coefficients of each feature in the regression model demonstrating the relative contribution of each feature to fitness residual (p value for regression coefficient of mRNA level = $3.5 \times 10^{-11}$, initiation rate = $2.5 \times 10^{-12}$, $TE_{GFP\ protein/mRNA}$ = $2.7 \times 10^{-9}$, codon decoding speed = $8.7 \times 10^{-3}$, mRNA folding energy = $1.5 \times 10^{-50}$, aSD velocity = $8.7 \times 10^{-3}$, hydrophobicity < $10^{-200}$, and amino acid synthesis cost = $5.4 \times 10^{-80}$). The sign of the contribution of each coefficient shows whether a feature is associated positively or negatively with fitness residuals. Error bars represent standard error of the coefficient estimation.

(C) Predicted fitness residuals of E. coli genes according to the regression model are correlated with their expression levels (Pearson correlation, r = 0.25, p = $2 \times 10^{-53}$), suggesting that natural selection shapes 5′ gene architectures in order to minimize costs of gene expression.

(D) Distribution of fitness residual scores for E. coli genes as predicted by regression model that was trained on either experimental or mock data. The experimentally based model predicts a significant, higher range of fitness residuals (p < $10^{-5}$), suggesting that the mechanisms that we elucidate with the synthetic library also apply on natural genes.

(E) Predicted fitness residuals of B. subtilis genes according to the regression model are correlated with their expression levels (Pearson correlation, r = 0.33, p = $10^{-93}$), suggesting that our model also applies for other bacteria species.

(F) Same as (D), only for B. subtilis genes.

cells in a dosage-dependent manner (Geiler-Samerotte et al., 2011). It is interesting to postulate that hydrophobic residues that promote aggregation can reduce the portion of properly folded, functional protein. Such futile protein synthesis might need to be compensated for by further costly production in order to reach the needed functional level of a certain protein.

We further demonstrate that there are sufficient degrees of freedom for a gene to evolve a cost-reducing architecture, even when its amino acid sequence is constant. Hence, our study suggests design elements that could be utilized both for better heterologous gene expression and by natural selection for the optimization of natural genes.

As such, our observations are also relevant to biotechnology and synthetic biology. Many times in such non-natural systems, there is a need to express a foreign gene, whose expression could deprive resources from the hosting cell. Our results allow the design of an optimized nucleotide sequence version for heterologous expression that minimizes the cost of production and, by that, reduces the burden on the cell while not compromising expression level.

**EXPERIMENTAL PROCEDURES**

See Supplemental Experimental Procedures for full description.

**CellPress**

## Library Architecture

The synthetic library was provided to us by Goodman et al. (2013) and is fully described there. In short, each variant in the library harbors a unique 5′ gene architecture that is composed of a promoter, a ribosome binding site, and an N′ terminus amino acid fusion of 11 amino acids followed by a super-folder GFP (sfGFP) gene. The library as a whole includes two promoters with either high or low transcription rates; three synthetic RBSs with strong, medium, or low translation initiation rates, as well as 137 different genomic RBSs that were defined as the 20 bp upstream to the ORF of 137 *E. coli* genes; and, finally, 137 coding sequences (CDSs) consisting of the first 11 amino acids from the same genes. Each CDS appears in the library in 13 different nucleotide sequences representing alternative synonymous forms. All combinations amounted in 14,234 distinct library variants.

## Competition Assay

Competition experiment was carried out by serial dilution. The library was grown on 1.2 mL of Lysogeny broth (LB) and 50 μg/mL kanamycin at 30°C, the exact same conditions that were used in Goodman et al. (2013) to measure GFP expression level. We grew six parallel, independent lineages, and each was diluted daily by a factor of 1:120 into fresh media (resulting in ~6.9 generations per dilution). This procedure was repeated for 12 days, and samples were taken from each lineage every 4 days (~27 generations), mixed with glycerol, and kept at −80°C.

## Fitness and Fitness Residual Estimations

Fitness effect is derived from the following equation:

$$f(t) = f(anc) \cdot (1+s)^t \approx f(anc) \cdot e^{st}$$

where $f$ is the variant frequency, $t$ is the generation number, and $s$ is the fitness effect.

To extract fitness effect, we took two independent approaches. First, we took the logarithm of the ratio between the frequency of a variant at a certain time point and its frequency at time zero. We then divided this value by the number of generations. This calculation was performed for both generation ~84 and generation ~56. See Supplemental Experimental Procedures for description of fitness calculation based on maximum likelihood. The two fitness-estimation methods were highly correlated (Figures S5A and S5B; r = 0.99, p < 10$^{-200}$) and resulted in the same conclusions throughout our analyses.

We then defined "fitness residual" of a variant as the difference between the observed fitness by FitSeq and the fitness predicted by a linear model given the variant's GFP expression level (see Supplemental Experimental Procedures for further details).

## Model for Estimating Translation Velocity Based on Anti-Shine Dalgarno Affinity

The Shine-Dalgarno affinity was calculated identically to Li et al. (2012). In short, for each position, we calculated the affinity of 8–11 bp upstream of that position (the distance between the ribosome A site and the aSD site) to the anti-Shine Dalgarno motif. The free energy of interaction between the aSD motif and the mRNA sequence (ΔG) was calculated for all possible 10-mer sequences for that position using the RNA annealing function from the ViennaRNA package algorithm (Lorenz et al., 2011), and the highest affinity (lowest energy) score was used. We calculated the affinity for all positions for which the annealing with the aSD motif resides in the 11 amino acid fusion (positions 19–33) and then transformed all affinities of a given variable sequence to estimated ribosomal velocity, as follows.

We converted the ΔG estimates into the equilibrium constant of the interaction, K, which represents the equilibrium between association ($k_f$) and dissociation ($k_b$). The elongation velocity ($v$) as the ribosome moves from current site $n$ to the n + 1 site is given by the harmonic mean of the dissociation reaction of site $n$ and the association reaction of site $n + 1$:

$$\frac{1}{v_{n \to n+1}} = \frac{1}{k_{b_n}} + \frac{1}{k_{f_{n+1}}} \qquad \text{Equation 1}$$

$$v_{n \to n+1} = \frac{k_{b_n} k_{f_{n+1}}}{k_{b_n} + k_{f_{n+1}}} \qquad \text{Equation 2}$$

We further assume that the association reaction rate is not dependent on the sequence, therefore, for every $n$, $k_{f_n} = k_f$, and that differences in affinity thus only reflect differences in dissociation constant displayed by various sequences. We then get a term for the ribosomal velocity at a specific position by the anti-Shine Dalgarno affinity:

$$v_{n \to n+1} = \frac{k_f \cdot k_f K^{-1}}{k_f(1+K^{-1})} = k_f \frac{e^{\frac{\Delta G}{RT}}}{1 + e^{\frac{\Delta G}{RT}}} \qquad \text{Equation 3}$$

To calculate the average ribosomal velocity across the entire N terminus fusion sequence of each library variant, we calculated the harmonic mean of the velocity values for all positions. See Supplemental Experimental Procedures for full description.

## ACCESSION NUMBERS

The accession number for all sequencing data reported in this paper is SRA: SRP092267.

## SUPPLEMENTAL INFORMATION

Supplemental Information includes Supplemental Experimental Procedures, five figures, and four tables and can be found with this article online at http://dx.doi.org/10.1016/j.molcel.2016.11.007.

## AUTHOR CONTRIBUTIONS

## ACKNOWLEDGMENTS

## REFERENCES

Akashi, H., and Gojobori, T. (2002). Metabolic efficiency and amino acid composition in the proteomes of Escherichia coli and Bacillus subtilis. Proc. Natl. Acad. Sci. USA *99*, 3695–3700.

Artieri, C.G., and Fraser, H.B. (2014). Accounting for biases in riboprofiling data indicates a major role for proline in stalling translation. Genome Res. *24*, 2011–2021.

Bennett, B.D., Kimball, E.H., Gao, M., Osterhout, R., Van Dien, S.J., and Rabinowitz, J.D. (2009). Absolute metabolite concentrations and implied enzyme active site occupancy in Escherichia coli. Nat. Chem. Biol. *5*, 593–599.

Bentley, W.E., Mirjalili, N., Andersen, D.C., Davis, R.H., and Kompala, D.S. (1990). Plasmid-encoded protein: the principal factor in the "metabolic burden" associated with recombinant bacteria. Biotechnol. Bioeng. 35, 668–681.

Bienick, M.S., Young, K.W., Klesmith, J.R., Detwiler, E.E., Tomek, K.J., and Whitehead, T.A. (2014). The interrelationship between promoter strength, gene expression, and growth rate. PLoS ONE 9, e109105.

Charneski, C.A., and Hurst, L.D. (2013). Positively charged residues are the major determinants of ribosomal velocity. PLoS Biol. 11, e1001508.

Charneski, C.A., and Hurst, L.D. (2014). Positive charge loading at protein termini is due to membrane protein topology, not a translational ramp. Mol. Biol. Evol. 31, 70–84.

Dana, A., and Tuller, T. (2014). The effect of tRNA levels on decoding times of mRNA codons. Nucleic Acids Res. 42, 9171–9181.

Dekel, E., and Alon, U. (2005). Optimality and evolutionary tuning of the expression level of a protein. Nature 436, 588–592.

Dong, H., Nilsson, L., and Kurland, C.G. (1995). Gratuitous overexpression of genes in Escherichia coli leads to growth inhibition and ribosome destruction. J. Bacteriol. 177, 1497–1504.

Emilsson, V., and Kurland, C.G. (1990). Growth rate dependence of transfer RNA abundance in Escherichia coli. EMBO J. 9, 4359–4366.

Gasch, A.P., and Werner-Washburne, M. (2002). The genomics of yeast responses to environmental stress and starvation. Funct. Integr. Genomics 2, 181–192.

Geiler-Samerotte, K.A., Dion, M.F., Budnik, B.A., Wang, S.M., Hartl, D.L., and Drummond, D.A. (2011). Misfolded proteins impose a dosage-dependent fitness cost and trigger a cytosolic unfolded protein response in yeast. Proc. Natl. Acad. Sci. USA 108, 680–685.

Gingold, H., and Pilpel, Y. (2011). Determinants of translation efficiency and accuracy. Mol. Syst. Biol. 7, 481.

Glick, B.R. (1995). Metabolic load and heterologous gene expression. Biotechnol. Adv. 13, 247–261.

Goodarzi, H., Nguyen, H.C.B., Zhang, S., Dill, B.D., Molina, H., and Tavazoie, S.F. (2016). Modulated expression of specific tRNAs drives gene expression and cancer progression. Cell 165, 1416–1427.

Goodman, D.B., Church, G.M., and Kosuri, S. (2013). Causes and effects of N-terminal codon bias in bacterial genes. Science 342, 475–479.

Heyer, E.E., and Moore, M.J. (2016). Redefining the translational status of 80S monosomes. Cell 164, 757–769.

Higgs, P.G., and Ran, W. (2008). Coevolution of codon usage and tRNA genes leads to alternative stable states of biased codon usage. Mol. Biol. Evol. 25, 2279–2291.

Ibarra, R.U., Edwards, J.S., and Palsson, B.O. (2002). Escherichia coli K-12 undergoes adaptive evolution to achieve in silico predicted optimal growth. Nature 420, 186–189.

Ingolia, N.T., Ghaemmaghami, S., Newman, J.R.S., and Weissman, J.S. (2009). Genome-wide analysis in vivo of translation with nucleotide resolution using ribosome profiling. Science 324, 218–223.

Kafri, M., Metzl-Raz, E., Jona, G., and Barkai, N. (2016). The cost of protein production. Cell Rep. 14, 22–31.

Kudla, G., Murray, A.W., Tollervey, D., and Plotkin, J.B. (2009). Coding-sequence determinants of gene expression in Escherichia coli. Science 324, 255–258.

Kyte, J., and Doolittle, R.F. (1982). A simple method for displaying the hydropathic character of a protein. J. Mol. Biol. 157, 105–132.

Li, G.-W., Oh, E., and Weissman, J.S. (2012). The anti-Shine-Dalgarno sequence drives translational pausing and codon choice in bacteria. Nature 484, 538–541.

Lorenz, R., Bernhart, S.H., Höner Zu Siederdissen, C., Tafer, H., Flamm, C., Stadler, P.F., and Hofacker, I.L. (2011). ViennaRNA Package 2.0. Algorithms Mol. Biol. 6, 26.

Marr, A.G. (1991). Growth rate of Escherichia coli. Microbiol. Rev. 55, 316–333.

Mohammad, F., Woolstenhulme, C.J., Green, R., and Buskirk, A.R. (2016). Clarifying the translational pausing landscape in bacteria by ribosome profiling. Cell Rep. 14, 686–694.

Plotkin, J.B., and Kudla, G. (2011). Synonymous but not the same: the causes and consequences of codon bias. Nat. Rev. Genet. 12, 32–42.

Qian, W., Yang, J.R., Pearson, N.M., Maclean, C., and Zhang, J. (2012). Balanced codon usage optimizes eukaryotic translational efficiency. PLoS Genet. 8, e1002603.

Rang, C., Galen, J.E., Kaper, J.B., and Chao, L. (2003). Fitness cost of the green fluorescent protein in gastrointestinal bacteria. Can. J. Microbiol. 49, 531–537.

Salis, H.M. (2011). The ribosome binding site calculator. Methods Enzymol. 498, 19–42.

Scott, M., Gunderson, C.W., Mateescu, E.M., Zhang, Z., and Hwa, T. (2010). Interdependence of cell growth and gene expression: origins and consequences. Science 330, 1099–1102.

Shah, P., and Gilchrist, M.A. (2011). Explaining complex codon usage patterns with selection for translational efficiency, mutation bias, and genetic drift. Proc. Natl. Acad. Sci. USA 108, 10231–10236.

Shah, P., Ding, Y., Niemczyk, M., Kudla, G., and Plotkin, J.B. (2013). Rate-limiting steps in yeast protein translation. Cell 153, 1589–1601.

Sharp, P.M., Tuohy, T.M., and Mosurski, K.R. (1986). Codon usage in yeast: cluster analysis clearly differentiates highly and lowly expressed genes. Nucleic Acids Res. 14, 5125–5143.

Subramaniam, A.R., Pan, T., and Cluzel, P. (2013). Environmental perturbations lift the degeneracy of the genetic code to regulate protein levels in bacteria. Proc. Natl. Acad. Sci. USA 110, 2419–2424.

Taniguchi, Y., Choi, P.J., Li, G.-W., Chen, H., Babu, M., Hearn, J., Emili, A., and Xie, X.S. (2010). Quantifying E. coli proteome and transcriptome with single-molecule sensitivity in single cells. Science 329, 533–538.

Tholstrup, J., Oddershede, L.B., and Sørensen, M.A. (2012). mRNA pseudo-knot structures can act as ribosomal roadblocks. Nucleic Acids Res. 40, 303–313.

Tuller, T., and Zur, H. (2015). Multiple roles of the coding sequence 5′ end in gene expression regulation. Nucleic Acids Res. 43, 13–28.

Tuller, T., Carmi, A., Vestsigian, K., Navon, S., Dorfan, Y., Zaborske, J., Pan, T., Dahan, O., Furman, I., and Pilpel, Y. (2010a). An evolutionarily conserved mechanism for controlling the efficiency of protein translation. Cell 141, 344–354.

Tuller, T., Waldman, Y.Y., Kupiec, M., and Ruppin, E. (2010b). Translation efficiency is determined by both codon bias and folding energy. Proc. Natl. Acad. Sci. USA 107, 3645–3650.

Tuller, T., Veksler-Lublinsky, I., Gazit, N., Kupiec, M., Ruppin, E., and Ziv-Ukelson, M. (2011). Composite effects of gene determinants on the translation speed and density of ribosomes. Genome Biol. 12, R110.

Vind, J., Sørensen, M.A., Rasmussen, M.D., and Pedersen, S. (1993). Synthesis of proteins in Escherichia coli is limited by the concentration of free ribosomes. Expression from reporter genes does not always reflect functional mRNA levels. J. Mol. Biol. 231, 678–688.

Wagner, A. (2005). Energy constraints on the evolution of gene expression. Mol. Biol. Evol. 22, 1365–1374.

Weinberg, D.E., Shah, P., Eichhorn, S.W., Hussmann, J.A., Plotkin, J.B., and Bartel, D.P. (2016). Improved ribosome footprint and mRNA measurements provide insights into dynamics and regulation of yeast translation. Cell Rep. 14, 1787–1799.

Wen, J.-D., Lancaster, L., Hodges, C., Zeri, A.-C., Yoshimura, S.H., Noller, H.F., Bustamante, C., and Tinoco, I. (2008). Following translation by single ribosomes one codon at a time. Nature 452, 598–603.

Yona, A.H., Bloom-Ackermann, Z., Frumkin, I., Hanson-Smith, V., Charpak-Amikam, Y., Feng, Q., Boeke, J.D., Dahan, O., and Pilpel, Y. (2013). tRNA genes rapidly change in evolution to meet novel translational demands. eLife 2, e01339.
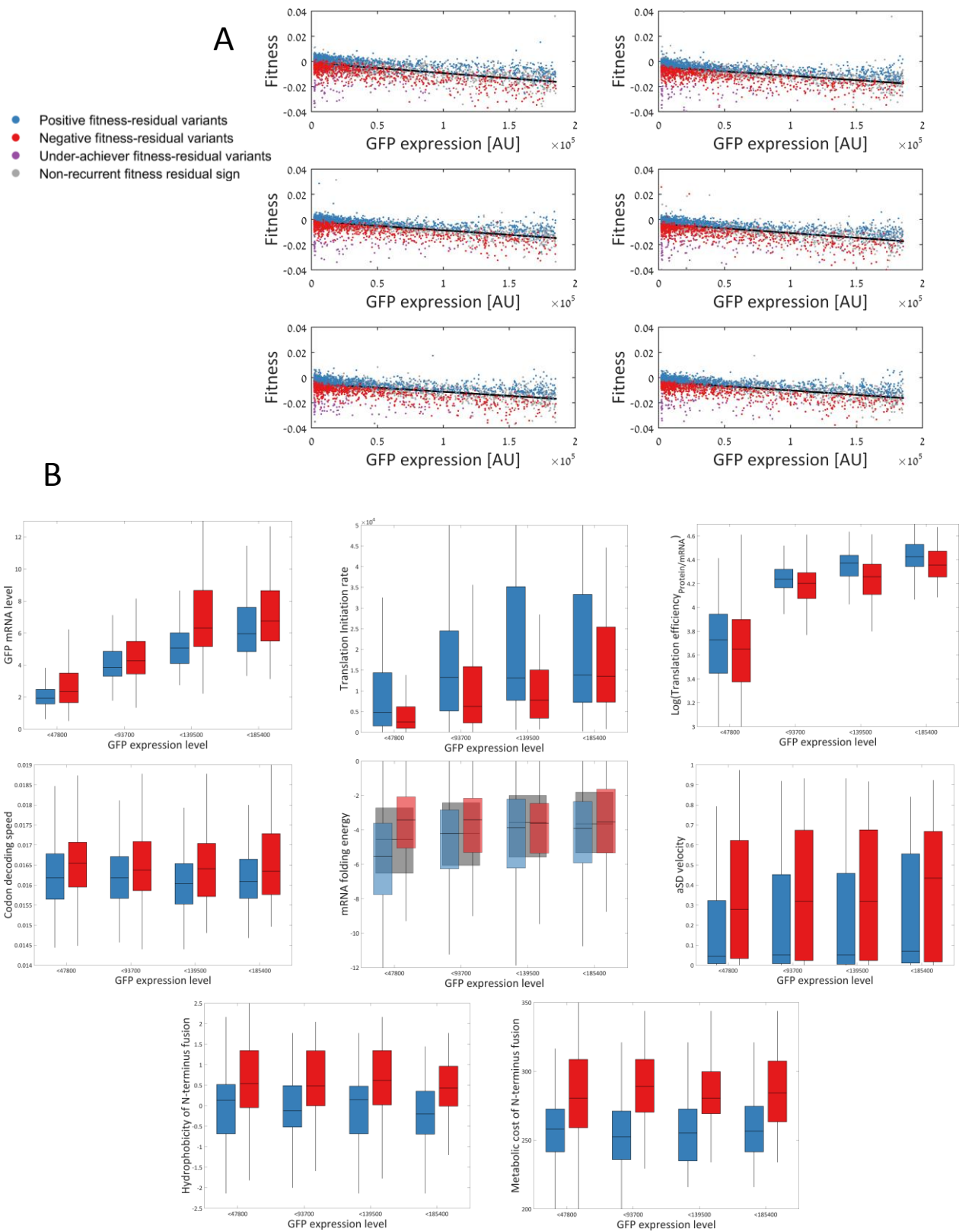
**Supplemental Information**

**Gene Architectures that Minimize**

**Cost of Gene Expression**

Idan Frumkin, Dvir Schirman, Aviv Rotman, Fangfei Li, Liron Zahavi, Ernest Mordret, Omer Asraf, Song Wu, Sasha F. Levy, and Yitzhak Pilpel
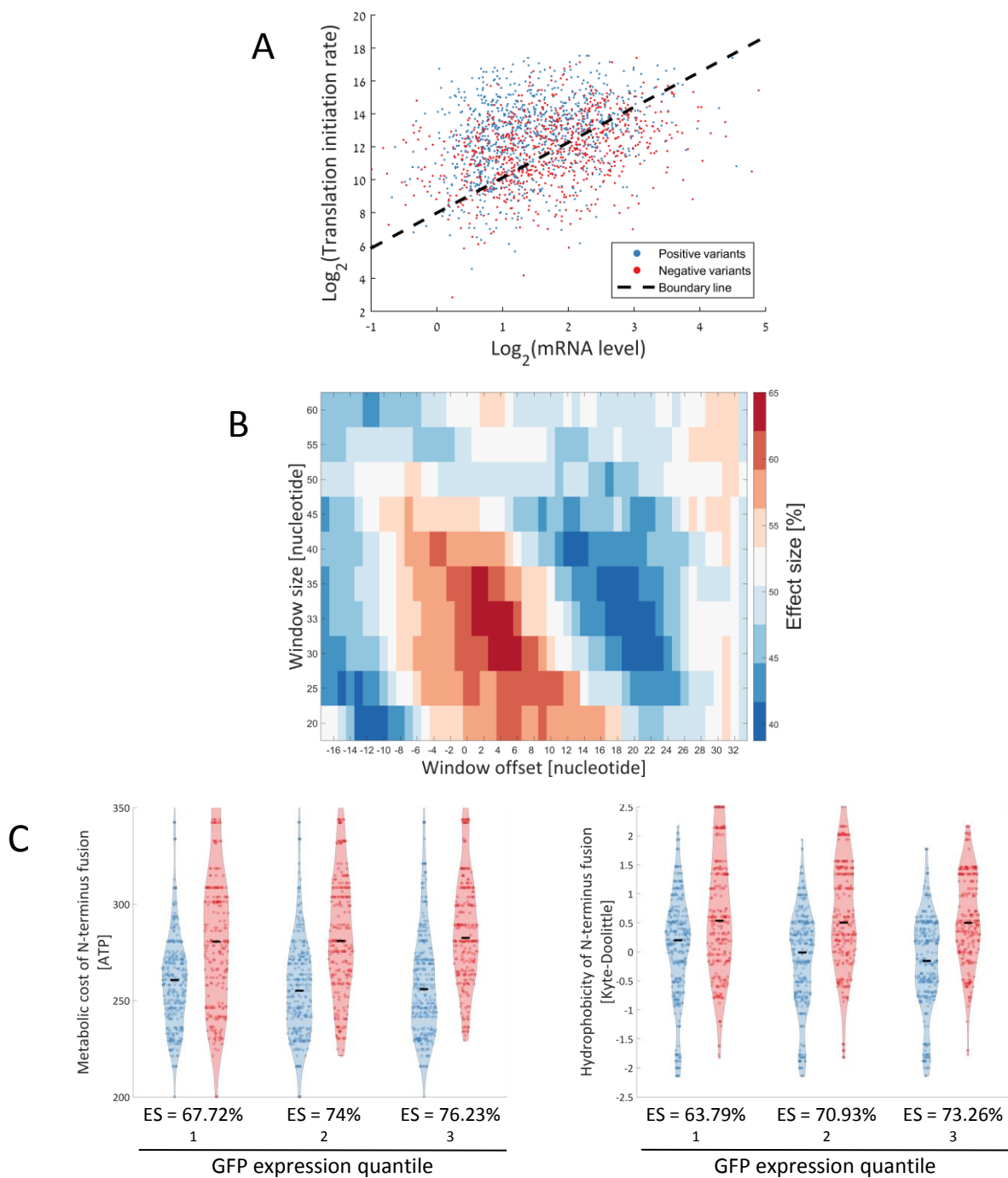
# Supplementary Figure 1

**Supplementary Figure 1, related to Figures 1-4**

**A| Correlation between fitness and GFP level for each FitSeq repeat.** The correlation between fitness and GFP expression level is presented for each independent competition of the synthetic library. The Pearson's r for each repeat is: -0.76, -0.76, -0.79, -0.78, -0.74 and -0.77, respectively. All p-Values are lower than $10^{-200}$.

**B| Contribution of each feature to fitness residual in bins of GFP expression level.** Library variants were binned according to GFP protein expression level and further split between positive and negative fitness residual variants. Positive (blue) and negative (red) variants from each GFP expression bin were then compared according to each of the eight features that affect fitness residual. Gray boxplots represent entire library variants.

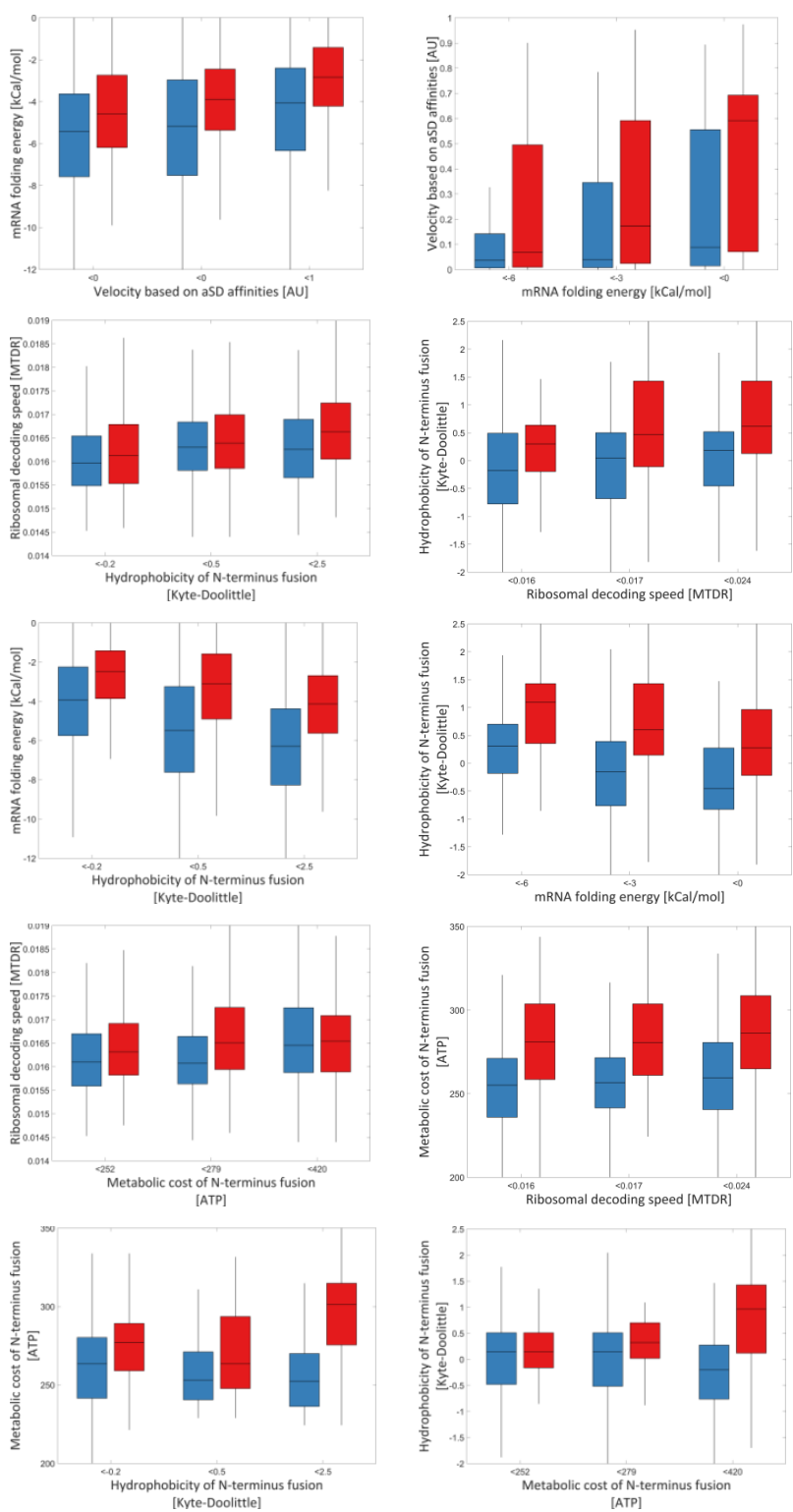**Supplementary Figure 2, related to Figure 2+3+4**

**A| Positive fitness residual variants have lower GFP mRNA levels but higher initiation rates compared to negative variants.** Positive (blues dots) and negative (red dots) variants are drawn according to their mRNA level of the GFP gene (X-axis) and their translation initiation rate (Y-axis). The dashed line represents the optimal linear boundary line between the positive and negative variants, as was computed by training an SVM classifier with a linear kernel function on all the variants, using the two axes as features.

**B| Color map of effect size for mRNA folding energy comparison between positive and negative fitness residual variants.** X-axis is the window starting position and Y-axis is the window size. The largest effect size of the difference in folding energy between positive and negative variants is observed when the window is positioned exactly at the variable region, just after the AUG codon.

**C| The higher the GFP expression, the more beneficial it is to utilize cheap or hydrophilic amino acids.**
**Left** Positive and negative variants were split into three equally-populated quantiles according to GFP expression levels. Then, the effect size for hydrophobicity between positive and negative variants was calculated for all quantiles. The difference in effect sizes between $1^{st}$ and $2^{nd}$ and between $1^{st}$ and $3^{rd}$ quantiles was calculated and found to be significant compared to random split of the variants (p-Values=0.018 and 0.0022, respectively). "ES" denotes Effect Size. **Right** Same as left, only for amino acid synthesis cost (p-Values=0.0088 and 0.0008).
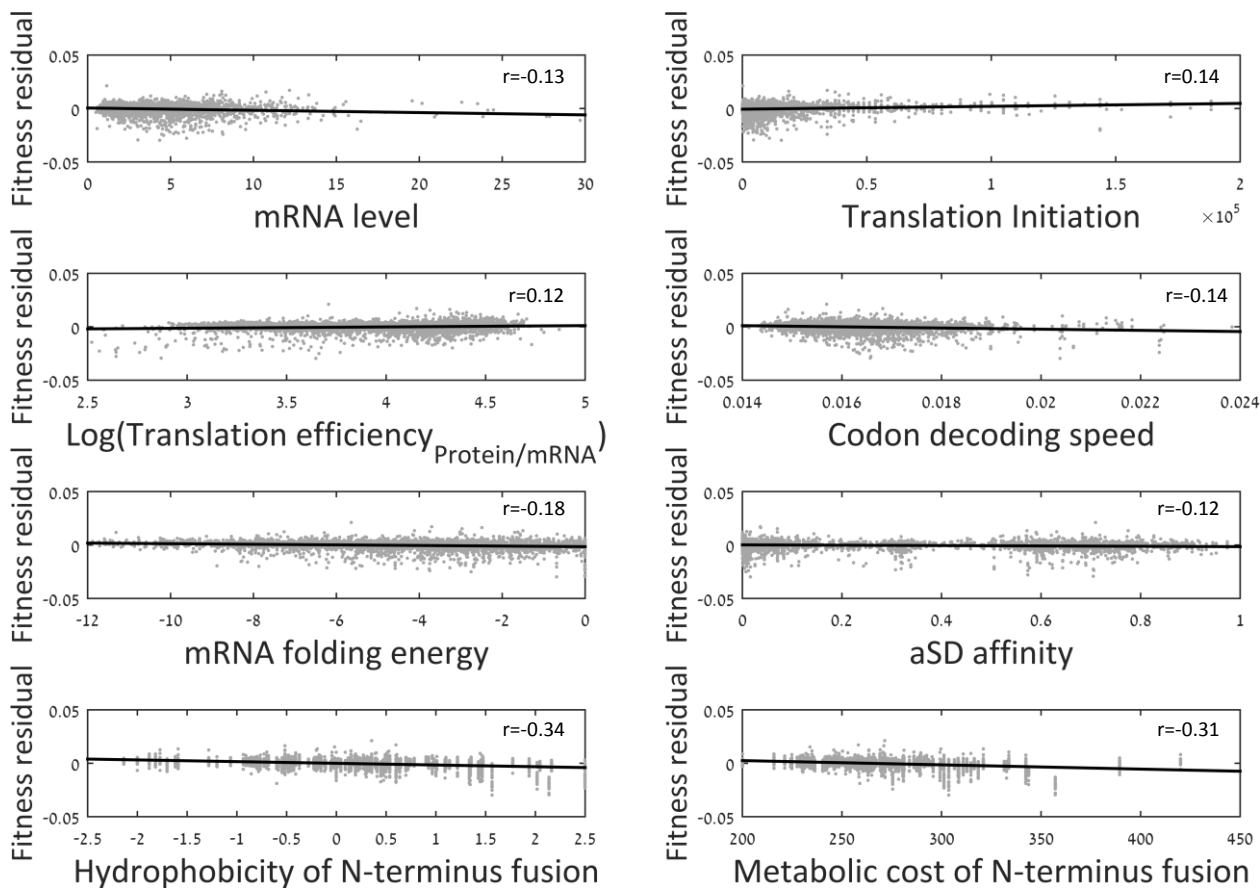
**Supplementary Figure 3, related to Figure 5 – Controlling the correlation between each feature**
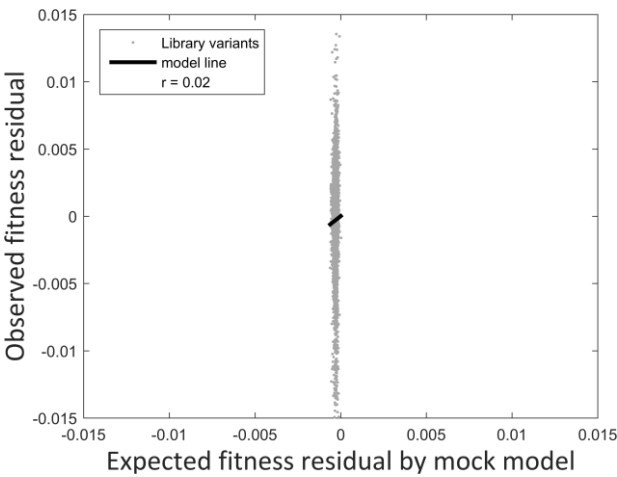To test the independent contribution of each feature to fitness residual, we checked the Pearson correlation between each pair of features (Figure 5). For features with r>0.1, we binned library variants according to the first feature, and then compared in each bin the difference between positive and negative fitness residual variants in the second feature. We then reciprocally binned the data according to the second feature and compared the first feature. List of feature pairs for which this analysis was performed: mRNA folding energy with aSD velocity, ribosomal decoding speed with hydrophobicity, mRNA folding energy with hydrophobicity, ribosomal decoding speed with cost of N-terminus fusion and hydrophobicity with cost of N-terminus fusion. For example, binned positive and negative fitness residual variants according to anti-Shine Dalgarno affinities still demonstrate mRNA folding energy difference in all bins (Wilcoxon rank-sum p-Value=$2.8 \cdot 10^{-5}$, $2 \cdot 10^{-7}$ and $2.5 \cdot 10^{-10}$). The reverse (bin according to folding energy and significant difference in aSD affinities) is also observed (Wilcoxon rank-sum p-Value=$8.9 \cdot 10^{-3}$, $4.5 \cdot 10^{-6}$ and $3 \cdot 10^{-11}$).

A



B



**Supplementary Figure 4, related to Figure 7**
**A| Individual sequence features demonstrate weak correlations with fitness residual**
Correlation of each feature, which was shown to differentiate between positive and
negative variants, with fitness residual. Pearson's r was calculated for each correlation.
**B| Mock model and cross validation for linear regression model**
Mock regression model that was trained on randomly scrambled library data fails to
accurately predict fitness residual.

**Supplementary Figure 5, related to Methods**
**A+B| comparison of fitness evaluation**
**A|** Comparison of fitness estimation of all library variants between the different repeats of the experiment.
**B|** The fitness estimations for all library variants with the two methods that were used in this study are highly correlated (see Methods, Pearson's r=0.99, p-Value<$10^{-200}$).
**C| Two methods of evaluation GFP expression level from FACS data**
In *Goodman et al.*, cells were sorted into 12 expression bins using FACS, and in each bin the relative frequency of each variant was measured using deep-sequencing. The estimated expression level of each variant was then calculated by computing the weighted geometric mean of the bins' median expression level, using the relative frequency of each variant at each bin as the bin's weight. In order to validate these data, we estimated the GFP expression level from the raw data by fitting gamma distribution parameters to the histogram of each variant's frequencies in all bins (see Methods). These two estimation methods are highly correlated (r=0.94, p-Value<$10^{-200}$), yet ~600 variants showed high expression levels according to the gamma fit method, while coming from the entire range of expression level using the geometric mean method. We thus excluded these variants from our analyses since their GFP measurement is not accurate.

**Table S1, related to Figure 1 –** Raw sequencing data.

**Table S2, related to Figure 1 –** Fitness per library variant.

**Table S3, related to Figures 2-4 –** Fitness residual and feature values per library variant.

**Table S4, related to Figure 7 –** Predicted fitness residual and feature values per *E. coli* or *B. subtilis* gene.

## Supplementary Materials and Methods File

### Library architecture

The synthetic library was provided to us by Goodman *et al.* (Goodman et al., 2013) and is fully described there. In short, each variant in the library harbors a unique 5' gene architecture that is composed of a promoter, a Ribosome Binding Sites (RBS) and an N'-terminus amino acid fusion of 11 amino acids followed by a sfGFP gene. The library as a whole includes: two promoters with either high or low transcription rate. Three synthetic RBSs with strong, medium, or low translation initiation rates, as well as 137 different genomic RBSs that were defined as the 20bps upstream to the ORF of 137 *E. coli* genes. And finally, 137 coding sequences (CDS) consisting of the first 11 amino acids from the same genes. Each CDS appears in the library in 13 different nucleotide sequences representing alternative synonymous forms. All combinations amounted in 14,234 distinct library variants.

### Competition Assay

Competition experiment was carried out by serial dilution. The library was grown on 1.2ml of LB + 50µg/ml kanamycin at 30°C, the exact same conditions as was used in Goodman *et al.* to measure GFP expression level. We grew six parallel, independent lineages and each was diluted daily by a factor of 1:120 into fresh media (resulting in ~6.9 generations per dilution). This procedure was repeated for 12 days and samples were taken from each lineage every four days (~27 generations), mixed with glycerol and kept at -80°C.

### Library preparation and sequencing

Plasmids from time zero (library "ancestor") and all other samples were purified with a QIAgene mini-prep kit and used as templates for PCR to amplify specifically the variable region of all variants in the population. To minimize PCR and sampling biases, we used a large amount of template, ~500ng of DNA, and a relatively short PCR of 26 rounds. The forward primer (sequence: CAGCTCTTCGCCTTTACGCATATG) was paired with 5 different reverse primers that

are one bp shifted from each other to insure that library complexity was high enough for Illumina sequencing:

R1: GACAATGAAAAGCTTAGTCATGGCG ; R2: ACAATGAAAAGCTTAGTCATGGCG

R3: CAATGAAAAGCTTAGTCATGGCG ; R4: AATGAAAAGCTTAGTCATGGCG

R5: ATGAAAAGCTTAGTCATGGCG

PCR products were then run on BluePipin to capture the correct amplicon size of ~140 bps and remove any un-specific amplicons. Then, DNA buffer was exchanged using Agencourt AmPure SPRI bead cleanup protocol. Hiseq library was prepared next using the sequencing library module from *Blecher-Gonen. et al.* 2013 (Blecher-Gonen et al., 2013). In short, blunt ends were repaired, Adenine bases were added to the 3' end of the fragments, barcode adapters containing a T overhang were ligated, and finally the adapted fragments were amplified. The process was repeated for each sample with a different Illumina DNA barcode for multiplexing, and then all samples were pooled in equal amounts and sequenced. We performed a 125bp paired end high output run on HiSeq 2500 PE Cluster Kit v4. Base calling is performed by RTA v. 1.18.64, and de-multiplexing is carried out with Casava v. 1.8.2, outputting results in FASTQ format.


**Data processing**

De-multiplexed data was received in the form of FASTQ files split into samples. First, SeqPrep (https://github.com/jstjohn/SeqPrep) was used to merge paired reads into a single contig, to increase sequence fidelity over regions of dual coverage. The size of each contig was then compared to the theoretical combined length of the forward primer, the reverse primer and the variable region of the variants. Next, the forward and reverse primers were found on each contig (allowing for 2 mismatches) and trimmed out. This step was performed for both the forward and reverse complement sequences of the contig, to account for non-directional ligation of the adaptors during library preparation. Then, the reverse primer was searched at the last 5 nucleotides of the contig to account for different primer lengths. After primers were

trimmed, the contig was tested again for its length to ensure no indels had occurred. Contigs were then compared sequentially to the entire library, comparing the sequence of each contig to the sequence of each variant. Any contig without a matching variant within two mismatches or less was discarded. Contigs with more than a single matching variant with the same reliability were also discarded due to ambiguity. Each contig that passed these filters was counted in key-value data structure, storing all variants in the library and their frequency in each sample. These data were then used for all downstream analyses. See raw data in Table S1, related to Figure 1.

**Fitness estimation**

Fitness effect is derived from the following equation:

$$f(t) = f(anc) \cdot (1 + s)^t \approx f(anc) \cdot e^{st}$$

Where f is the variant frequency, t is the generation number and $s$ is the fitness effect.

To extract fitness effect, we took two independent approaches. First, we took the logarithm of the ratio between the frequency of a variant at a certain time point and its frequency at time zero. We then divided this value by the number of generations. This calculation was performed both for generation ~84 and generation ~56. See fitness per variant in Table S2, related to Figure 1.

Second, we derived fitness by employing a Maximum-Likelihood (ML) algorithm on all frequency measurements along the competition experiment per variant. A key challenge to accurately estimating each variant's fitness over many generations is that the mean fitness of the population (against which each variant competes) changes (improves) over time. This is caused because more fit variants expand in the population at the expense of other (less fit) strains. To overcome this challenge, we use a Poisson likelihood maximization strategy (see full description below). Briefly, we make a first fitness estimate of each variant using a simple log-linear regression over the first three time points. Based on these estimations, the initial relative frequencies of each variant, and a noise model that accounts for experimental errors (Levy et al., 2015), we estimate the expected trajectory of each variant and compare this to the measured trajectory. We next make small changes to our fitness estimates, repeat this

comparison, accept updated fitness estimates if they better fit the data (higher likelihood), and perform this procedure iteratively until fitness estimates are stable (maximized likelihood).

These two fitness-estimation methods were highly correlated (Supplementary Figure 5A+B, r=0.99, p-Value<$10^{-200}$) and resulted with the same conclusions throughout our analyses.

**GFP expression level estimation**

GFP expression levels were taken from *Goodman et al.* (Goodman et al., 2013) data, in which it was calculated using the method described in *Kosuri et al.* (Kosuri et al., 2013). In short, cells were sorted into 12 expression bins using FACS, and in each bin the relative frequency of each variant was measured using deep-sequencing. The estimated expression level of each variant was then calculated by computing the weighted geometric mean of the bins' median expression level, using the relative frequency of each variant as the bin's weight.

In order to validate this data, we estimated the GFP expression level from the raw data in *Goodman et al.* by fitting gamma distribution parameters (suggested before as a model to capture noise, or spread of expression values of a gene within an isogenic population (Friedman et al., 2006)) to the histogram of each variant's frequencies in all bins. This gamma distribution follows this equation: $P(x) = \frac{x^{a-1}e^{-\frac{x}{b}}}{b^{a}\Gamma(a)}$ where $\Gamma$ denotes the gamma function.

These two estimation methods are highly correlated (Supplementary Figure 5C, r=0.94, p-Value<$10^{-200}$). However, we noticed that ~600 variants showed high expression levels according to the gamma fit method, while coming from the entire range of expression level using the geometric mean method. When closely examining these cases, we noticed that the source for the disagreement between the two methods is that these variants were observed only in two bins, with one of them being the highest bin, and the other not being the second highest. Therefore, we decided that the expression estimation for these variants is unreliable and excluded them from our analyses.

**Calculation of fitness residuals and classifying variants according to their positive or negative fitness residual sign**

We defined "fitness residual" of a variant as the difference between the observed fitness by FitSeq and the fitness predicted by a linear model given the variant's GFP expression level. To calculate fitness residual, we performed the following steps:

First, we filtered out variants that demonstrate a lower GFP level than $2^{11}$[AU], since below this threshold the GFP measurement method is not sensitive to accurately measure GFP. We also excluded variants with a GFP level above $2^{17.5}$[AU], as above this threshold the measurement method saturates. Notably, only variants with the "high promoter" were included in the analysis, since almost all "low promoter" variants did not pass the protein level filter. This decision was essential as the few "low promoter" variants that did pass this threshold show biased values of sequence features, such as a very low GC-content, which could mask real signals.

Next, we fitted a linear regression model between fitness and GFP expression levels for each of the six independent FitSeq repeats separately at each of the last two time points (generations ~56 and ~84). Then, variants for which fitness residual was in the top or bottom 5% were excluded and a new regression line was fitted in order to reach a better fit. These outliers were excluded only for the sake of fitting a regression line and were still included in our downstream analyses. Then, a fitness residual score was calculated for each variant at each repeat of the experiment and on each of the two time points.

We then split the variants into two groups: positive or negative fitness residual variants. To account for random processes (experimental errors and drift), "positive" or "negative" class was assigned for a given variant only if it showed a positive/negative fitness residual sign in at least 5 lineages in both time points. The set of all the above filters resulted in 975 variants in the positive variant group and 815 in the negative variant group.

Since we noticed that some of the negative variants have extreme negative fitness residual values, we further classified them as "underachievers". Underachiever variants were defined as variants that repeatedly showed fitness residual values in the bottom 5% of the entire library.

Similar to the positive/negative classification, a variant is assigned as "underachiever" only if it is found in the bottom 5% in at least 5 out of the 6 linages in both time points, which resulted in 80 variants.

**Parameter comparison between two fitness residual groups**

A one-sided Wilcoxon rank-sum test was used to compare the distributions of different sequence parameters between the positive and negative fitness residual groups. We also tested the effect size of each parameter using the "Probability of superiority" method (Ruscio, 2008) that calculates the probability to randomly choose a member from group A with a higher value than a random member from group B.

To compare between effect sizes according to GFP expression levels, we split the positive and negative variant groups into three quantiles according to GFP expression levels. Then, the effect size for hydrophobicity or amino acid synthesis cost between positive and negative variants were calculated for each quantile. We then performed an empirical p-Value estimation by randomly choosing three data sets with the same number of variants, and computed the effect size at each set. This sampling was performed $10^4$ times, and p-Value was estimated by counting the number of times the difference in effect sizes between the first and second sets and between the first and third sets were lower in the real data than the difference in effect sizes of the random groups.

**Calculating translation initiation rate per variant**

We estimated the translation initiation rate of each variant with the "RBS calculator" (Espah Borujeni et al., 2014; Salis et al., 2009), which simulates initiation rates given a UTR and a coding sequence. This calculation is achieved by using a bio-mechanic model combining the affinity to the anti-Shine Dalgarno sequence of the ribosome, mRNA secondary structure of the UTR and coding sequence, and steric interference of the ribosome and the mRNA.

**Mean of the Typical Decoding Rates (MTDR) estimation**

To evaluate codon-decoding times by the ribosome we used a published index of Mean of the Typical Decoding Rates (MTDR) values (Dana and Tuller, 2014), which were derived from ribosome profiling data (Li et al., 2012). MTDR values for each of the 61 sense codons are driven from measured ribosome density, when the ribosome A site is on a codon, averaged over all the appearances of the codon within mRNAs. This measurement estimates the translation speed of each codon, and it correlates significantly (r=0.46 for *E. coli)* with tRNA availability. The final score given to each variant was the harmonic mean of its MTDR values of the first 11 amino acids.

**Folding energy estimation of mRNA secondary structure**

We calculated folding energy of mRNA secondary structure for each variant by using the ViennaRNA package algorithm (Lorenz et al., 2011). Each sequence was computed by a sliding window, whose starting position ranged from position -18 to position 32. The calculation was repeated with different window sizes (20-60bps). All calculations were done assuming a temperature of $30^{o}$C.

**Model for estimating translation velocity based on anti-Shine Dalgarno affinity**

The Shine-Dalgarno affinity was calculated identically to Li *et al.* (Akashi and Gojobori, 2002). In short, for each position we calculated the affinity of 8-11bps upstream of that position (the distance between the ribosome A site and the aSD site) to the anti-Shine Dalgarno motif. The free energy of interaction between the aSD motif and the mRNA sequence ($\Delta$G) was calculated for all possible 10mer sequences for that position using the RNA annealing function from the ViennaRNA package algorithm (Lorenz et al., 2011), and the highest affinity (lowest energy) score was used. We calculated the affinity for all positions for which the annealing with the aSD motif resides in the 11-amino acid fusion (positions 19-33) and then transformed all affinities of a given variable sequence to estimated ribosomal velocity as follows.

We converted the ΔG estimates into the equilibrium constant of the interaction, K by:

(i) $\qquad K = e^{-\frac{\Delta G}{RT}}$

Where $\Delta G$ denotes the SD affinity, $R$ denotes the gas constant and $T$ denotes the temperature.

This equilibrium constant, at the n$^{th}$ codon along a sequence, is defined in turn, given the association reaction rate ($k_f$) which represents the association to the current site ($n$), and a dissociation reaction ($k_b$) that represents the dissociation to the current site as:

(ii) $\qquad K = \frac{k_{f_n}}{k_{b_n}}$

The elongation velocity ($v$) as the ribosome moves from current site $n$ to the n+1 site is given by the harmonic mean of the dissociation reaction of site $n$ and the association reaction of site $n+1$:

(iii) $\qquad \frac{1}{v_{n \to n+1}} = \frac{1}{k_{b_n}} + \frac{1}{k_{f_{n+1}}}$

(iv) $\qquad v_{n \to n+1} = \frac{k_{b_n} k_{f_{n+1}}}{k_{b_n} + k_{f_{n+1}}}$

We further assume that the association reaction rate is not dependent on the sequence, therefore for every $n$, $k_{f_n} = k_f$. Introducing equations (i)-(ii) to the equation (iv), results in a term for the ribosomal velocity at a specific position by the anti-Shine Dalgarno affinity:

(v) $\qquad v_{n \to n+1} = \frac{k_f \cdot k_f K^{-1}}{k_f (1 + K^{-1})} = k_f \frac{e^{\frac{\Delta G}{RT}}}{1 + e^{\frac{\Delta G}{RT}}}$

To calculate the average ribosomal velocity across the entire N-terminus fusion sequence of each library variant, we calculated the harmonic mean of the velocity values for all positions. The analysis was performed also at a codon resolution, taking into account only positions of the sequence that are the first nucleotide of codons, which yielded similar results to the nucleotide-based analysis.

**Amino acid property estimation of N-terminus fusion amino acids**

Hydrophobicity of each 11-amino acid N-terminus peptide was calculated according to its score on the Kyte-Doolittle scale (Kyte and Doolittle, 1982). Amino acid cost was derived from Akashi and Gojobori (Akashi and Gojobori, 2002) in the form of the amount of energy consumed for its production in high energy ATP or GTP bonds. Cost was either evaluated per amino acid or summed for the whole peptide.

Supply of amino acids were derived from *Bennet et al.* (Bennett et al., 2009), which measured cellular concentrations of amino acid in exponnentially grown *E. coli*. Notably, two amino acids are missing from this table (Gly & Cys), and two amino acids are indistinguishable (Lys & Ile). Therefore, those 4 amino acids were excluded from the this analysis. The demand per amino acid was calculated by multiplying the frequency of each amino acid in each *E. coli* gene by the median ribosome profiling score of the gene (Li et al., 2012). The sum of all genes was defined as the total amino acid demand.


**Amino acid enrichment in positive and negative variant groups**

To calculate the frequency of the various amino acids in the collective proteome in either the positive or the negative fitness residual group, we counted the occurrences of each amino acid in each variant. We then summed this number for each amino acid across all variants in each group and divided the sum by the number of variants in each group multiplied by 11. To quantify the relationship between amino acid enrichment and energetic-cost or availability we calculated the frequency ratio of each amino acid by dividing the amino acid frequency of the positive fitness residual group by the frequency of the negative group. We then calculated the Pearson correlation between the log2 amino acid enrichment ratio and the amino acid energetic-cost or their availability.

**Comparing fitness residual among variants with the same N-terminus fusion by Δfitness-residual**

We defined Δfitness-residual as the difference between the fitness residual of a given variant with the average fitness residual of the variant group with the same N-terminus amino acid fusion. Therefore, Δfitness-residual measures the expression cost of a variant normalized to its GFP expression level and its N-terminus amino acid sequence. We then spilt each variant group of the same N-terminus fusion to above-average and below-average variants and calculated for each sub group the mean value for six features (RNA levels, translation initiation rates, translation efficiency, codon decoding speed (MTDR), mRNA secondary structure, and anti-Shine Dalgarno affinity). For each feature, the mean value of the below-average (x-axis) and above-average (y-axis) Δfitness-residual groups were depicted as a scatter plot, in which each point represents a different N-terminus fusion. Then, the deviance of all dots from the identity (X=Y) line was calculated and tested for significance with a one-tailed Student's t-test. To compute an effect size for this enrichment, we used Cohen's d: $d = \frac{\bar{x}_{high} - \bar{x}_{low}}{S}$ where $\bar{x}_{high \backslash low}$ represents the mean of the feature in the above- or below-average group, and $S$ represents the standard deviation of the feature in the entire set of library variants that was used in this study.

**A multiple linear regression model to predict fitness residual**

We performed a multiple linear regression using all eight features as independent variables (RNA levels, translation initiation rate, translation efficiency (GFP protein/mRNA), mRNA secondary structure, codon-decoding speed, aSD affinities, amino acid metabolic cost and hydrophobicity) and the mean fitness residual across six repeats of FitSeq as the dependent variable. The regression yielded a coefficient for each feature, which were all used in order to predict fitness residual of a given variant.

As a negative control for this model we randomly shuffled each of the features in the library, trained a mock model on this shuffled library, and computed the Pearson correlation coefficient between the observed fitness residual and the expected fitness residual according to the mock model. In order to compute a p-Value we repeated this process $10^5$ times, and counted the

number of times the correlation coefficient from the mock model was higher than the correlation coefficient from the real model.

To predict fitness residual of natural *E. coli* and *B. subtilis* genes, a second regression model was performed, in which we excluded translation efficiency (due to lack of data for the entire ~4000 *E. coli* genes) and hydrophobicity (due to the fact that hydrophobic motifs in membrane proteins are functional, hence including this feature might lead to wrong estimation of membrane proteins). We also used Lasso regularization and feature selection method (Tibshirani, 1996) with Matlab's "lassoglm" function from the "Statistics and Machine Learning" toolbox to avoid overfitting of the model. The $\lambda$ value was chosen as the value for which the deviance was one standard deviation higher than the minimum deviance achieved in a 1000-fold cross validation. Out of the six features used for this model, none were excluded by Lasso method. This model performed well in predicting fitness residual of the library variants and a cross validation test resulted in correlation of r=0.3 (p-Value=$10^{-10}$).

This model was then used to predict fitness residual scores for natural *E. coli* (strain MG1655) and *B. subtilis* (strain 168) genes. For each gene of these species, we computed a score for each feature of the model. We used RNA levels for *E. coli* from a previous RNA-seq experiment in which cells ware grown in LB and were harvested at the logarithmic growth phase. We used published RNA data for *B. subtilis* (Cohen et al., 2016). Translation initiation rates was computed with the same initiation rate model as was used for the library variants (Espah Borujeni et al., 2014; Salis et al., 2009). mRNA secondary structure, codon-decoding speed and aSD affinities were calculated as explained for the library variants. MTDR values for both species were taken from published data (Dana and Tuller, 2014). Amino acid metabolic cost was calculated as the mean value for the entire protein, and for both species the same cost was assigned for each amino acid (Akashi and Gojobori, 2002). Protein expression levels for both species were taken from the integrated datasets in Pax-Db (Wang et al., 2012). See Table S4, related to Figure 7 for feature and fitness residual scores for genes of these two organisms.

As a negative control for the prediction of fitness residual for natural *E. coli* genes, we generated a mock model by randomly shuffling each of the features in the library, training a linear regression model on this shuffled library and using it to predict fitness residual for all *E.*

*coli* genes. We then compared the standard deviation of the fitness residual predictions by the real model to the one of the mock model. This analysis was repeated $10^5$ times to compute a p-Value for the chance of the real model to show a higher standard deviation than the mock model.

**Cross validation sets**

Cross validation tests of the regression model were performed by randomly choosing training and test sets, in proportions of 70% and 30% of the entire library variants, respectively. In order to account for the fact that some of the information lays in the amino acid sequence, the training/test sets were also separated by the N-terminus amino acid peptide sequences with 41 peptide sequences (~30%) chosen as test set, and the rest as training sets. 10-fold cross validation was performed by randomly generating ten different pairs of training and test sets. The results are based on the average across these 10 repeats.

**RNA fitness residual calculation**

To evaluate mRNA fitness residuals we repeated the same calculation as described for fitness residual only with the mRNA levels instead of protein levels placed on the x-axis.

**Fitness residual and feature values per variant**

See Table S3, related to Figures 2-4.

**Bibliography**

Akashi, H., Gojobori, T., 2002. Metabolic efficiency and amino acid composition in the proteomes of Escherichia coli and Bacillus subtilis. Proc. Natl. Acad. Sci. U. S. A. 99, 3695–3700. doi:10.1073/pnas.062526999

Bennett, B.D., Kimball, E.H., Gao, M., Osterhout, R., Van Dien, S.J., Rabinowitz, J.D., 2009. Absolute metabolite concentrations and implied enzyme active site occupancy in Escherichia coli. Nat. Chem. Biol. 5, 593–9. doi:10.1038/nchembio.186

Blecher-Gonen, R., Barnett-Itzhaki, Z., Jaitin, D., Amann-Zalcenstein, D., Lara-Astiaso, D., Amit, I., 2013. High-throughput chromatin immunoprecipitation for genome-wide mapping of in vivo protein-DNA interactions and epigenomic states. Nat. Protoc. 8, 539–54. doi:10.1038/nprot.2013.023

Cohen, O., Doron, S., Wurtzel, O., Dar, D., Edelheit, S., Karunker, I., Mick, E., Sorek, R., 2016. Comparative transcriptomics across the prokaryotic tree of life. Nucleic Acids Res. gkw394. doi:10.1093/nar/gkw394

Dana, A., Tuller, T., 2014. The effect of tRNA levels on decoding times of mRNA codons. Nucleic Acids Res. 1–11. doi:10.1093/nar/gku646

Espah Borujeni, A., Channarasappa, A.S., Salis, H.M., 2014. Translation rate is controlled by coupled trade-offs between site accessibility, selective RNA unfolding and sliding at upstream standby sites. Nucleic Acids Res. 42, 2646–2659. doi:10.1093/nar/gkt1139

Friedman, N., Cai, L., Xie, X.S., 2006. Linking stochastic dynamics to population distribution: an analytical framework of gene expression. Phys. Rev. Lett. 97, 168302. doi:10.1103/PhysRevLett.97.168302

Goodman, D.B., Church, G.M., Kosuri, S., 2013. Causes and effects of N-terminal codon bias in bacterial genes. Science 342, 475–9. doi:10.1126/science.1241934

Kosuri, S., Goodman, D.B., Cambray, G., Mutalik, V.K., Gao, Y., Arkin, A.P., Endy, D., Church, G.M., 2013. Composability of regulatory sequences controlling transcription and translation in Escherichia coli. Proc. Natl. Acad. Sci. U. S. A. 110, 14024–9. doi:10.1073/pnas.1301301110

Kyte, J., Doolittle, R.F., 1982. A simple method for displaying the hydropathic character of a protein. J. Mol. Biol. 157, 105–32.

Levy, S.F., Blundell, J.R., Venkataram, S., Petrov, D. a., Fisher, D.S., Sherlock, G., 2015. Quantitative evolutionary dynamics using high-resolution lineage tracking. Nature advance on. doi:10.1038/nature14279

Li, G.-W., Oh, E., Weissman, J.S., 2012. The anti-Shine-Dalgarno sequence drives translational pausing and codon choice in bacteria. Nature 484, 538–41. doi:10.1038/nature10965

Lorenz, R., Bernhart, S.H., Höner Zu Siederdissen, C., Tafer, H., Flamm, C., Stadler, P.F., Hofacker, I.L., 2011. ViennaRNA Package 2.0. Algorithms Mol. Biol. 6, 26. doi:10.1186/1748-7188-6-26

Ruscio, J., 2008. A probability-based measure of effect size: robustness to base rates and other factors. Psychol. Methods 13, 19–30. doi:10.1037/1082-989X.13.1.19

Salis, H.M., Mirsky, E. a, Voigt, C. a, 2009. Automated design of synthetic ribosome binding sites to control protein expression. Nat. Biotechnol. 27, 946–950. doi:10.1038/nbt.1568

Tibshirani, R., 1996. Regression Shrinkage and Selection via the Lasso. J. R. Stat. Soc. 58, 267–

288.

Wang, M., Weiss, M., Simonovic, M., Haertinger, G., Schrimpf, S.P., Hengartner, M.O., von Mering, C., 2012. PaxDb, a database of protein abundance averages across all three domains of life. Mol. Cell. Proteomics 11, 492–500. doi:10.1074/mcp.O111.014704

**Simulation for measurement errors in GFP expression level**

To calculate fitness residual, we correlated fitness to expression level in order to learn the expected fitness of each library variant. Variants were classified as positive/negative fitness residual variants, if their fitness value was repeatedly above/below the expected line, respectively. This approach allowed us to factor out the effect of GFP expression level on fitness and elucidate mechanisms that reduce cost at a given expression level. Yet, a potential bias in our method could arise because of experimental errors in the GFP measurement. Indeed, *Goodman et al.*(Goodman et al., 2013) and *Kosuri et al.*(Kosuri et al., 2013) report an estimated coefficient of variation $\left(\frac{\sigma}{\mu}\right)$ of 0.22 for the GFP level in each variant. We set to assess the potential effect of such measurement error on fitness residual estimation.

We simulated our experimental design with a range of measurement errors for GFP expression level (see below a description of the simulation). This simulation allowed us to evaluate how many variants would be wrongly classified with either positive or negative fitness residual simply due to error in GFP expression levels measurements.

Our results show that for all simulated error levels, even those that far exceed the actual reported error level, we observed much less positive and negative variants in the simulation compared with the actual group size of the positive and negative variants in the FitSeq experiment (Simulation Figure 1A). This result means that our classification of variants into positive and negative fitness residual groups could not be due to error in measurements.

Additionally, our simulation predicts that GFP measurement error alone would result in a negative variant group that is larger than the positive variant group and only upon introduction of fitness residual signal to the simulation, more positive than negative variants were observed (Simulation Figure 1B). Reassuringly, our data resulted in a greater number of positive variants compared to negative, suggesting that library variants show a real phenomenon of fitness residual that minimize cost of gene expression.

Next, we turned to test whether the features that we discovered to differ between positive and negative variants, and thus affect fitness residual, could be observed

due to measurement errors (at the reported error level of *Goodman et al.*). For all features, except GFP mRNA levels (p-Value=0.8), we observed that the effect size that separates between the positive and negative variant groups is much higher than would appear simply because of experimental errors (Simulation Figure 1C). P-Values for initiation rate<$10^{-4}$, translation efficiency$_{protein\backslash mRNA}$<$10^{-4}$, codon decoding speed=$3.2x10^{-3}$, mRNA folding=$2x10^{-4}$, aSD velocity<$10^{-4}$, amino acid metabolic cost<$10^{-4}$ and hydrophobicity<$10^{-4}$. These results mean that the molecular mechanisms we revealed in this work are not observed due to an experimental error, but rather reflect a genuine biological phenomenon relating to expression cost.

Regarding mRNA levels, we present three arguments for its relevance to fitness residual: First, translation efficiency, defined as GFP protein/mRNA, at the variant level was still observed as significant in this analysis, a result which strengthens our claim that producing more proteins per mRNA reduces cost of gene expression. Second, we calculated "RNA fitness residual" based on fitness and two independent GFP mRNA measurements (Simulation Figure 1D and see methods) and observed that positive variants demonstrated lower mRNA levels compared to negative (Effect size=55.13%, rank-sum p-Value=$7.2x10^{-5}$), suggesting that higher mRNA levels are costly and reduce fitness residual. Third, we observed 80 variants with consistent extremely low fitness residual ("underachievers", see main text). While these very low fitness residual variants cannot be explained by measurement error (they do not appear in the simulation), they also demonstrate even higher mRNA levels than the negative variant group. All of these points suggest that mRNA level, as the other eight parameters, indeed reduce expression cost and increase fitness.

**Detailed description of the simulation**

For each single run of the simulation, 12 independent repeats are performed that simulate the 12 sampling points we used to classify each library variant with either positive or negative fitness residual. The simulation steps are as follows:

1. GFP expression level is randomly assigned for 4115 simulated variants from a log uniform distribution of GFP levels, which is similar to the distribution of the synthetic library we used in this study.

2. Fitness score is assigned to each simulated variant, according to the observed correlation between GFP expression level and fitness in our experiment.

   (i)    $\widehat{f_i} = a\widehat{GFP_i} + b$

$\widehat{f_i}$ – fitness predicted from GFP level according to linear model

$\widehat{GFP_i}$ – GFP expression levels drawn from the experimental distribution of GFP levels

a,b – confidents of the linear model as extracted from the measured data.

3. For each simulated variant, 12 independent measurement errors are added to the assigned fitness, in order to simulate the 12 repeats (6 from each time point) that were used for classifying the library variants. This fitness measurement error is drawn from a normal distribution with a mean of zero and a standard deviation of 0.03, N(0, 0.03). This SD was used as it is the mean SD we observed for the library variants based on the independent repeats of our experiment.

   (ii)    $f_i = \widehat{f_i} + N(0,0.03) = a\widehat{GFP_i} + b + N(0,0.03)$

$f_i$ - simulated fitness

4. A measurement error is added to the GFP level of each simulated variant by drawing a measurement error from a normal distribution. Since the absolute size of the measurement error is dependent on expression level (higher expressions mean larger errors) the simulated measurement error is chosen from a normal distribution with a mean of zero and an SD that is the multiplication of the simulated GFP level by the noise factor, $N_x$, which is a parameter of the simulation.

(iii) $\quad GFP_i = \widehat{GFP_i} + N\left(0, N_x \cdot \widehat{GFP_i}\right)$

$GFP_i$ – simulated GFP levels

$N_x$ – GFP noise factor

5. All simulated variants are then classified with positive or negative variants with the same approach as described in the methods section and the size of each group is counted.

6. The Pearson correlation between simulated GFP levels and simulated fitness scores is also calculated and recorded.

The above steps describe a single run of the simulation. We performed $10^4$ runs to calculate p-Value to the group size we observed in our study.

We then turned to simulate the effect size of each feature we observed to affect fitness residual:

For each simulated variant (as described in steps 1-6) we assigned a random feature score (taken from the actual values in the library) based on its randomly assigned GFP expression level. We performed this step while maintaining the correlation between expression level and the specific feature. For example, in the synthetic library we used in this study, mRNA folding energy is correlated with GFP expression with r=0.15. We maintained this correlation in the random assignment of folding energies to the simulated variants.

In order to produce a correlation $r$ between a given feature and the simulated GFP levels, we used the following method:

   a. Using Matlab's Statistics and Machine Learning toolbox copularnd() function we generated two random vectors with a normal distribution and 4115 samples each: $\{U_1, U_2\}$ with a correlation $r$ between them.

   b. Each of the vectors was sorted and a vector of the indices of $U$ mapping to the sorted vector was returned, such that: $U_i(I_i) = S_i$, where $S_i$ is the sorted vector, and $I_i$ are the indexes of $U_i$ ordered according to their rank.

c. The feature vector and the GFP vector were also sorted, we mark their sorted version as $S_{feature}$, $S_{GFP}$ respectively.

d. The correlated vectors of the feature and GFP vectors $X_{feature}$, $X_{GFP}$ were created by sampling the sorted vectors, using the indexes mapping to the sorted correlated normal vectors.
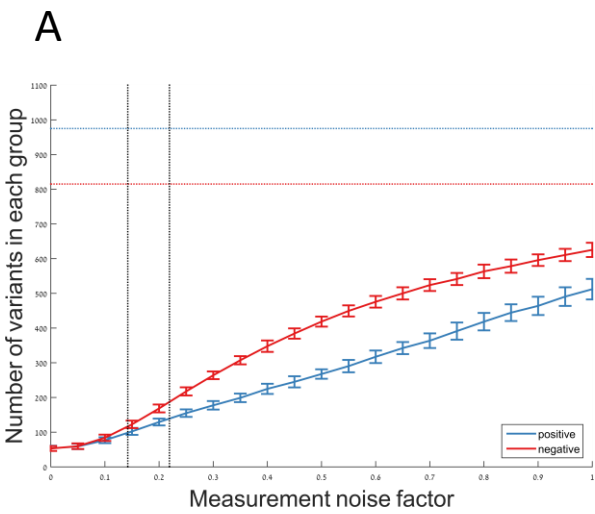
$$X_{feature} = S_{feature}(I_1)$$
$$X_{GFP} = S_{GFP}(I_2)$$

After creating a pair of vectors representing the feature values and simulated GFP values which have the same correlation as the measured feature have with the GFP expression levels, we repeated steps 2-5 in order to extract fitness residual values and positive/negative classification for the new permuted GFP vector.

We repeated the above process for each of the eight features.

7. Then for each feature, we calculated the effect size between the positive and negative fitness residual variants in the simulation. Since there are only experimental errors and no real signal in the simulation, this measured effect size is the threshold for our experimental design. To calculate p-Value to the effect size that we observed in the experiment data, we performed $10^4$ runs of the simulation and counted the number of times the effect size was higher than the one observed in the experiment.
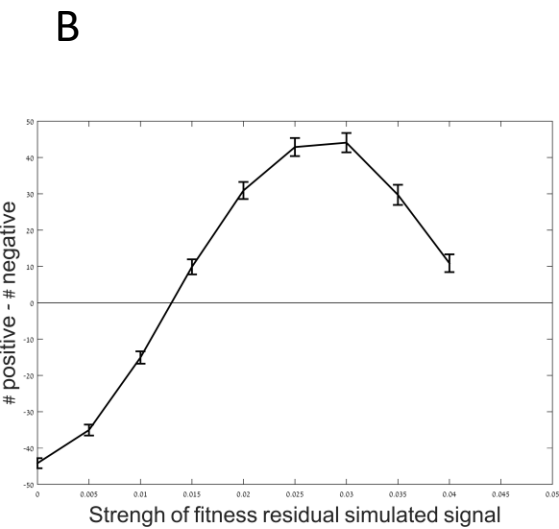
Simulation Figure 1

**A**

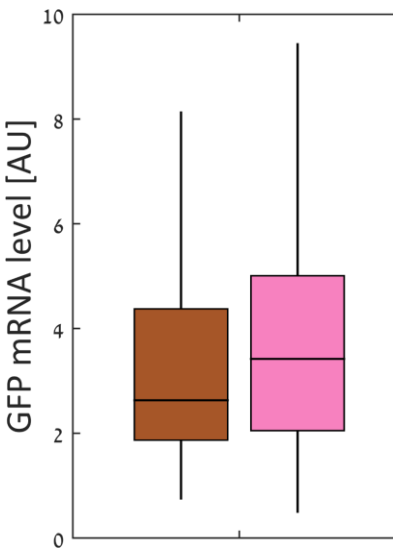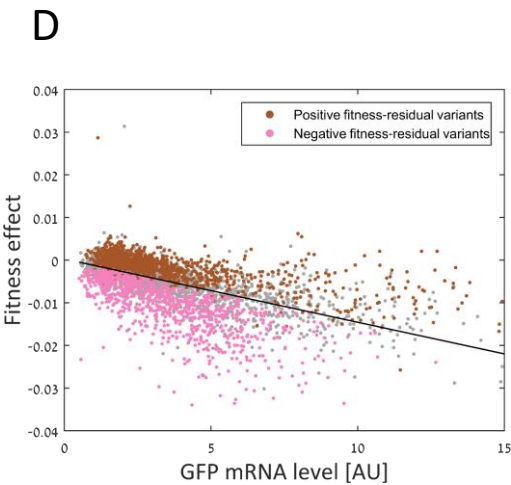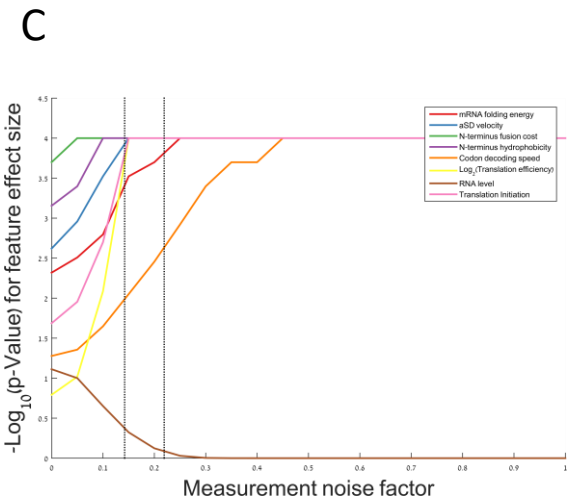Blue horizontal line - number of positive variants in the library.
Red horizontal line - number of negative variants in the library.
Left vertical line - noise factor corresponding to measured r between fitness and expression level in our data.
Right vertical line – noise factor as estimated from *Kosuri et al*.

**B**

Noise factor = 0.22 (as measured by from *Kosuri et al*.)

**C**

**D**

**Poisson likelihood maximization strategy to estimate fitness**

Here, each genetic design is called a lineage and the total number of lineages, $l$, is 14234. Let $n_i(t)$ be the cell number of lineage $i$ at the $t$-th generation, and $\overline{x}(t)$ be the mean fitness of the population at the $t$-th generation. In a limited-resource environment, the growth of lineage $i$ follows,

$$n_i(t) = n_i(0) \cdot e^{\int_0^t (x_i - \overline{x}(\tau)) d\tau},$$

where $\overline{x}(t) = \frac{\sum_{i=0}^l n_i(t) \cdot x_i}{\sum_{i=0}^l n_i(t)}$.

Therefore, the growth of lineage $i$ can be rewritten as,

$$n_i(t + \Delta t) = n_i(0) \cdot e^{\left(\int_0^t (x_i - \overline{x}(\tau)) d\tau + \int_t^{t+\Delta t} (x_i - \overline{x}(\tau)) d\tau\right)} \approx n_i(t) \cdot e^{\Delta t \cdot (x_i - \overline{x}(t))}.$$

We estimate the fitness of all lineages using the following method:

1. Let $x_i$ be fitness of lineage $i$. Let $r_i(t)$ and $f_i(t)$ be the experimental read number and read frequency of lineage $i$ at the $t$-th generation, respectively, i.e., $f_i(t) = \frac{r_i(t)}{\sum_{i=0}^l r_i(t)}$. Let $\Delta t$ be the number of generations passed between two bottlenecks. Here, $\Delta t \approx 28$. Make an initial guess of $x_i$ by linear regression of $\ln f_i(0)$, $\ln f_i(\Delta t)$, and $\ln f_i(2\Delta t)$. Note that we only do the linear regression for the lineage with an initial experimental read number larger than 10, i.e., $r_i(0) > 10$. Lineages with lower read numbers are too noisy to get an accurate estimation.

2. Let $\hat{n}_i(t)$ be the estimated cell number of lineage $i$ at the $t$-th generation. Assume that $\hat{n}_i(0) = \frac{r_i(0) \cdot N}{\sum_{i=0}^l r_i(0)}$, where $N$ is the total cell number after bottleneck. Here, $N = 9.37 \times 10^7$. Calculate $\hat{n}_i(t)$ and $\overline{x}(t)$ for the first 4 sequencing time points using,

$$\begin{cases} \overline{x}(k\Delta t) = \frac{\sum_{i=0}^l \hat{n}_i(k\Delta t) \cdot x_i}{\sum_{i=0}^l \hat{n}_i(k\Delta t)}, \\ \hat{n}_i(k\Delta t + \Delta t) = \hat{n}_i(k\Delta t) \cdot e^{\Delta t \cdot (x_i - \overline{x}(k\Delta t))}, \end{cases} \quad k = 0,1,2,3.$$

3. Let $\hat{r}_i(t)$ be the estimation of $r_i(t)$. Thus,

$$\hat{r}_i(k\Delta t) = \frac{\hat{n}_i(k\Delta t) \cdot \sum_{i=0}^l r_i(k\Delta t)}{\sum_{i=0}^l \hat{n}_i(k\Delta t)}, \quad k = 0,1,2,3.$$

4.  Define the Poisson likelihood function as,

$$F(x_1, x_2, \cdots, x_l) = \sum_{i=1}^{l} \sum_{k=0}^{3} \ln\left(\frac{\hat{r}_i(k\Delta t)^{r_i(k\Delta t)} \cdot e^{-\hat{r}_i(k\Delta t)}}{r_i(k\Delta t)!}\right).$$

Here, $\frac{\hat{r}_i(t)^{r_i(t)} \cdot e^{-\hat{r}_i(t)}}{r_i(t)!}$ gives the probability of observing the experimental read number $r_i(t)$ given the estimated read number $\hat{r}_i(t)$.

5.  Obtain the optimal fitness for all lineages by maximizing $F(x_1, x_2, \cdots, x_l)$. We use the *fminunc* function in the MATLAB Optimization Toolbox to do the optimization. The function *fminunc* uses Broyden–Fletcher–Goldfarb–Shanno (BFGS) algorithm with cubic line search to solve unconstrained nonlinear optimization problems. BFGS algorithm is an iterative method, which seeks a stationary point (with the derivative or gradient of the function being zero) of a function as Newton's method. The BFGS algorithm is a fast-converging algorithm and we find that all replicates nearly converge by 21 iterations (Maximum Likelihood Figure 1A).

This Poisson likelihood optimization method provides a more accurate estimation of fitness compared with log-linear regression method. We define the relative error of the estimation of read number of lineage $i$ as $\frac{1}{4} \cdot \sum_{k=0}^{3} \frac{2|\hat{r}_i(k\Delta t) - r_i(k\Delta t)|}{\hat{r}_i(k\Delta t) + r_i(k\Delta t)}$ and then calculate the relative error between the measured lineage trajectories and the lineage trajectories that would be expected based on our fitness estimates. We find that the fitnesses estimated using Poisson likelihood optimization have lower errors than the log-linear regression method (Maximum Likelihood Figure 1B).
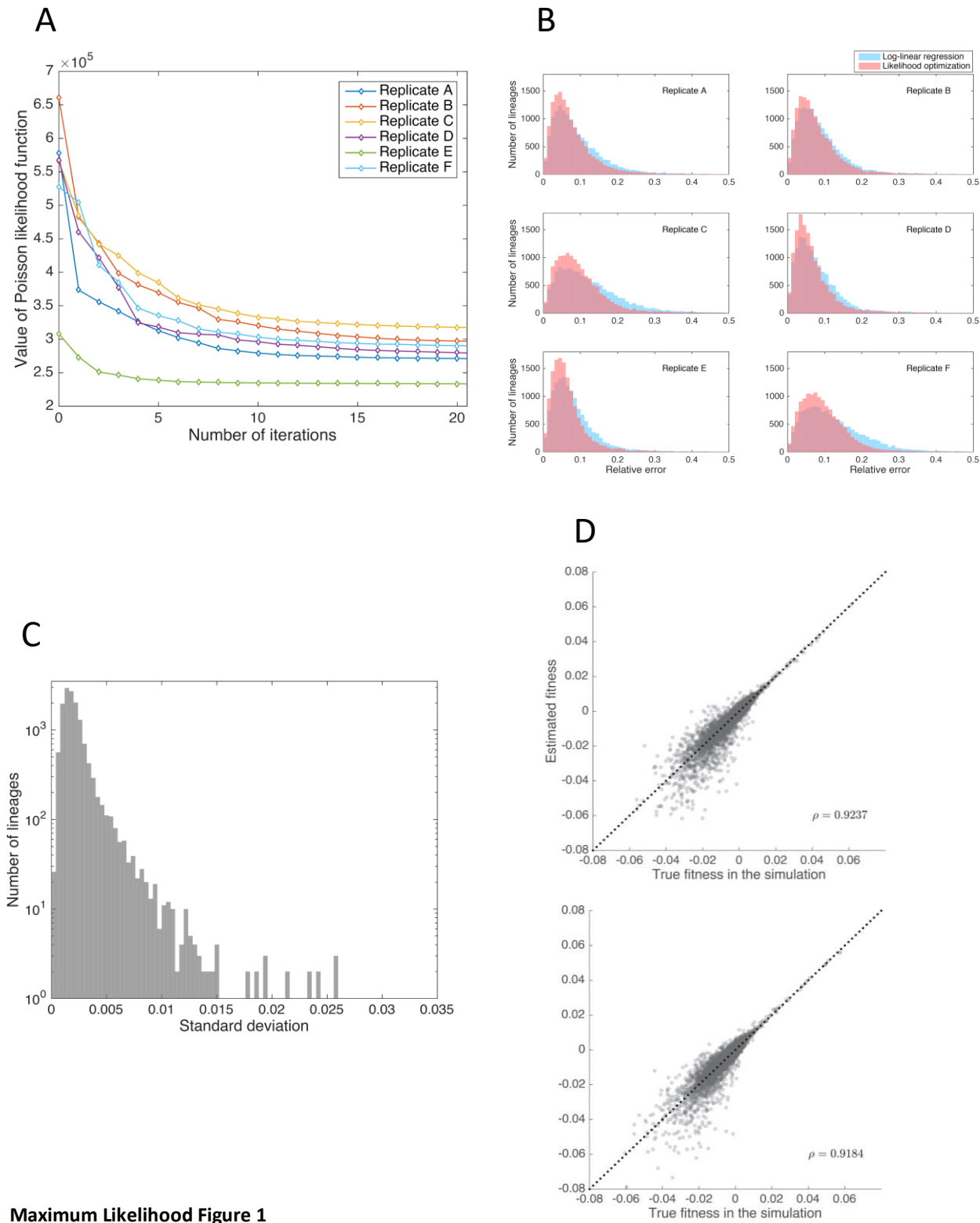
To test the consistency of the estimated fitnesses, we compared the Poisson likelihood optimization fitness estimates between all replicates. We find that replicates generally correlate well (Pearson's correlation > 0.76 for all replicates), with higher fitness designs having higher correlations between replicates. Additionally, we find that variance between fitness estimates across the six replicates is generally low (Maximum Likelihood Figure 1C).

To validate that the method worked as expected, we ran a simulation of the evolutionary process that would be expected in replicates E and F, given the fitness of each cell in lineage $i$ has the same fitness $x_i$ and that $x_i$ is the fitness estimated from real data using the Poisson likelihood optimization method. We assumed that there are no mutations and that the offspring per generation of an individual follows Poisson distribution with mean $1 + x_i$. Let $n_i^{simu}(t)$ be the cell number of lineage $i$ at the $t$-th generation in the simulation. In the

simulation, the cell number of lineage $i$ at the beginning is set as $n_i^{simu}(0) = \frac{r_i(0) \cdot N}{\sum_{i=0}^{l} r_i(0)}$, where $r_i(0)$ is the initial read number of lineage $i$. The evolution is simulated for 84 generations. To simulate the sequencing process where the data is sequenced every 28 generations, we let $r_i^{simu}(t)$ be the read number of lineage $i$ at the $t$-th generation, assuming that $r_i^{simu}(t)$ follows the Poisson distribution with mean $\frac{n_i^{simu}(t) \cdot r_i(t)}{\sum_{i=0}^{l} n_i^{simu}(t)}$, where $r_i(t)$ is the initial read number of lineage $i$. We then used the read number data obtained from these simulations to estimate the fitness for each lineage using the Poisson likelihood optimization method. We find that our fitness estimates correlate extremely well with the true (assigned) fitness of each lineage in both simulations (Maximum Likelihood Figure 1D).

**Maximum Likelihood Figure 1**
**A| Convergence of Poisson likelihood function.** A plot of the value of the Poisson likelihood function at each iteration for all six growth replicates.
**B| Distribution of relative error.** The distribution of the relative error for each replicate pooled growth using Poisson likelihood optimization method (pink) and the log-linear regression method (blue).
**C| Histogram of standard deviations of estimated fitness across six replicates using the Poisson likelihood optimization method.**
**D| Poisson likelihood optimization method performance on simulated data.** Scatter plots of the true and estimated fitnesses for evolutionary simulations with starting parameters that (up) match replicate E, or (down) match replicate F. $\rho$ is the Pearson correlation coefficient.