

## Estimating the errors in $r_{\text{cat}}$

### Flux Variability Analysis

One of the main concerns in our approach for calculating the maximal catalytic rate of an enzyme, is how to get a representative flux value to plug in the equation:

$$r_{\text{cat}} = \frac{v_i}{E_i} \quad (1)$$

Currently, we solve parsimonious FBA (pFBA), which minimizes the total sum of fluxes. It was shown that pFBA agrees with transcriptional data from *E. coli* in 97% of the time, and is considered a reasonable approach from flux predictions among the FBA community.

However, the flux values we get for each reaction are random sample from a possible flux distribution and may not represent the "real" flux through a given reaction. To reduce biases in our  $r_{\text{cat}}$  estimations, we can consider parsimonious *FVA*, which returns a range of fluxes for each reaction. It works in the following manner:

- (1) solve pFBA
- (2) set all reaction boundaries accordingly
- (3) for each reaction in the model:
  - (3a) minimize the reaction flux ( $v_{\text{min}}$ )
  - (3b) maximize the reaction flux ( $v_{\text{min}}$ )

The range of flux for reaction  $i$  is then:

$$v_{\text{max},i} - v_{\text{min},i} \quad (2)$$

$$(3)$$

Taking the mean of the  $v$  range may be a better approach for  $r_{\text{cat}}$  estimations. Error bars (or more precisely, confidence intervals) can be described as:

$$\overline{r_{\text{cat},i}} = \frac{v_{\text{max},i}}{E_i} \quad (4)$$

$$r_{\text{cat},i} = \frac{v_{\text{min},i}}{E_i} \quad (5)$$

$$(6)$$

Most of our reactions seem to have a very small dynamic range of fluxes, resulting in almost zero flux-associated errors (figure 1). Usually, if a reactions have small flux variability it means that they are biomass related, thus must support a given amount of flux to allow growth. Nevertheless, this is not the

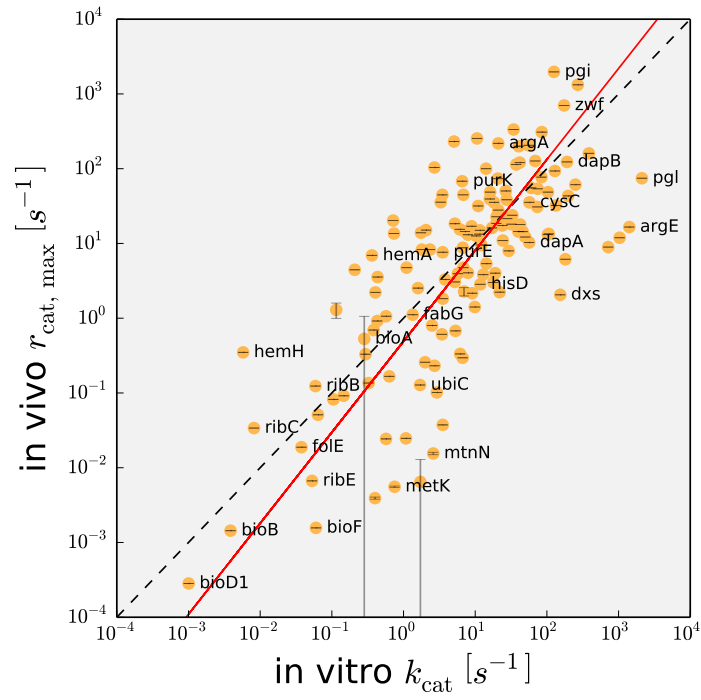


Figure 1:  $k_{\text{cat}}$  to  $r_{\text{cat}}^{\text{max}}$  correlation using  $v_{\text{mean}}$  as flux value for  $r_{\text{cat}}^{\text{max}}$  calculation. Dashed black line represents  $y = x$  and red line represents best fit by orthogonal regression. correlation:  $r^2 = 0.556$ ,  $p_{\text{val}} < 10^{-24}$ ,  $\text{rmse} = 0.6$

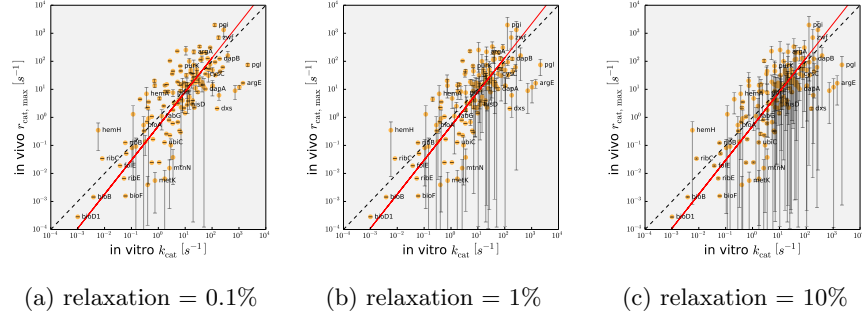


Figure 2

case for most of the 135 reactions we analyze (or at least I cannot see how they are biomass related). It is possible that the minimization of fluxes constrain, on top of the other constrains (growth rate and oxygen uptake rate) leave very little room for flux deviations. We can further relax the parsimonious constraint, that is, allow the sum of fluxes to be  $x\%$  larger. A common relaxing factor is in the range of 0.1-5%. Figure 2

### Standard deviation of the maximum

We can ask how reliable  $r_{cat}^{max}$  is. It may be that the maximal catalytic rate across all conditions. For example if the enzyme copy number in a particular condition was reported to be extremely small,  $r_{cat}^{max}$  will be disproportionately high. For this, we can calculate the standard deviation of the top three  $r_{cat}$  values:

$$\text{std}([r_{max}^{1st}, r_{max}^{2nd}, r_{max}^{3rd}]),$$

Figure 3a shows the resulting error bars. A reason for such large deviations is that the catalytic rate of many enzymes is log normal distributed, thus taking the std of the data results in huge errors. We can consider using:

$$\text{std}(\log([r_{max}^{1st}, r_{max}^{2nd}, r_{max}^{3rd}])),$$

yet then the errors are very small as we are only using three numbers (fig 3b).

### What Next?

What should we do? Maybe we can only show the positive error bars in figure 2, assuming that the minimal flux is not really biologically feasible, i.e., the algorithm may be using "irrelevant" metabolic pathways, allowing it to minimize the flux through the reactions significantly. It will look like this (fig 4):

