

Optimal Learning for Scientific Discovery

Peter Frazier, Jialei Wang, Pu Yang
School of Operations Research & Information Engineering
Cornell University

Mike Burkart*, Nathan Gianneschi*, Mike Gilson**, Nick Kosa*, Lori Tallorin*
Department of Chemistry & Biochemistry (*)
Skaggs School of Pharmacy and Pharmaceutical Sciences (**)
University of California San Diego

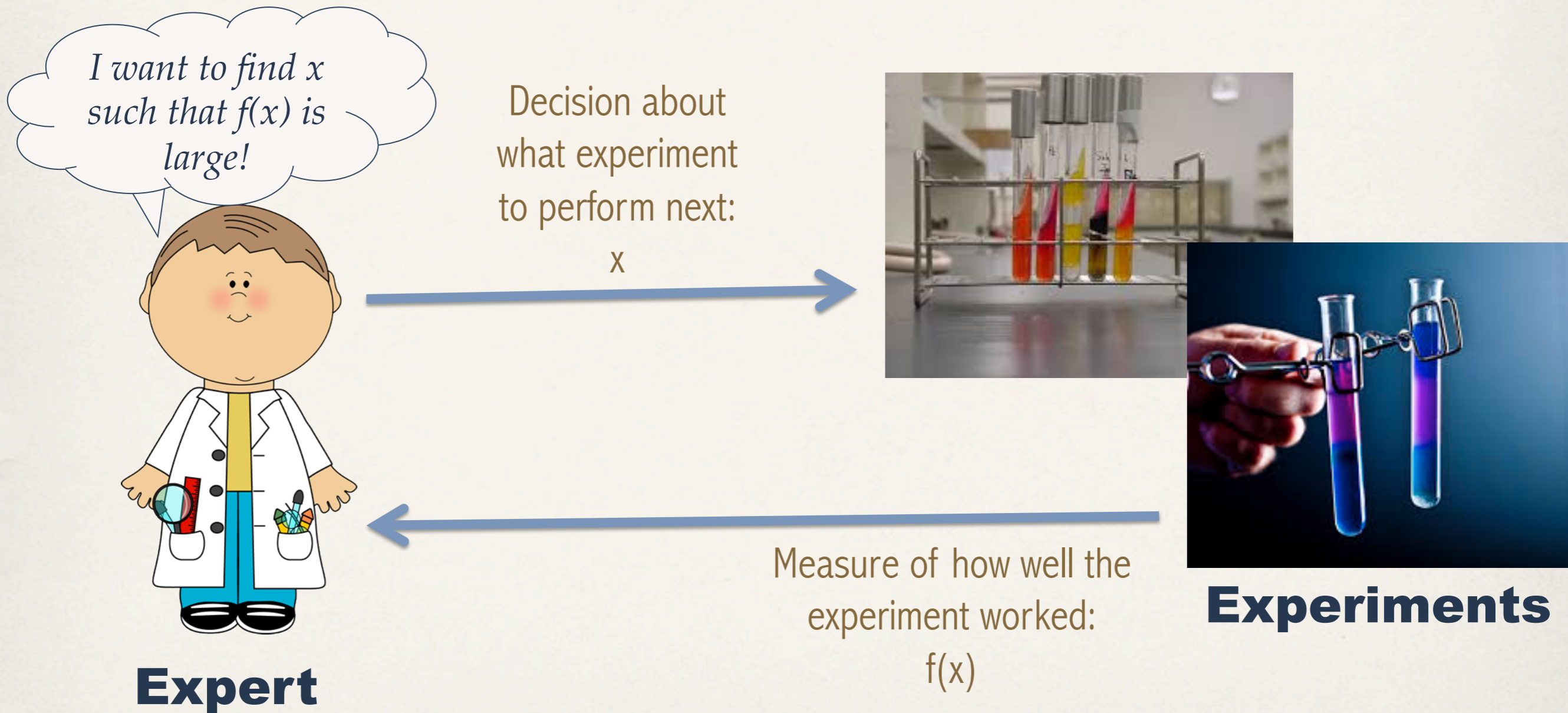
INFORMS Computing Society Conference, Richmond VA

January 11, 2015

Funded by AFOSR and NSF



Chemistry uses trial & error

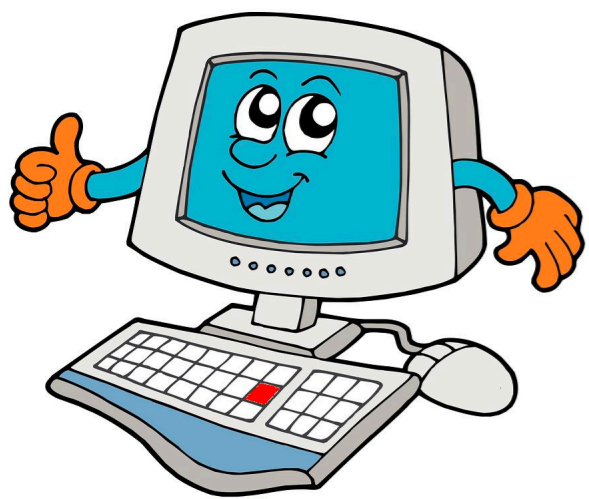


Trial & error is like optimization

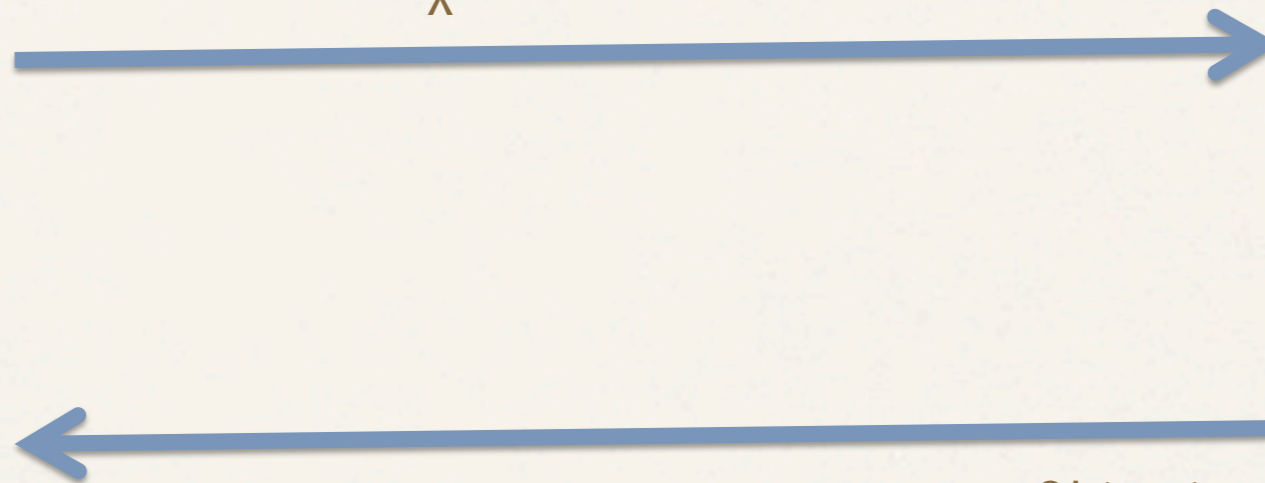
I want to find x such that $f(x)$ is large!

Decision about what point to evaluate next:

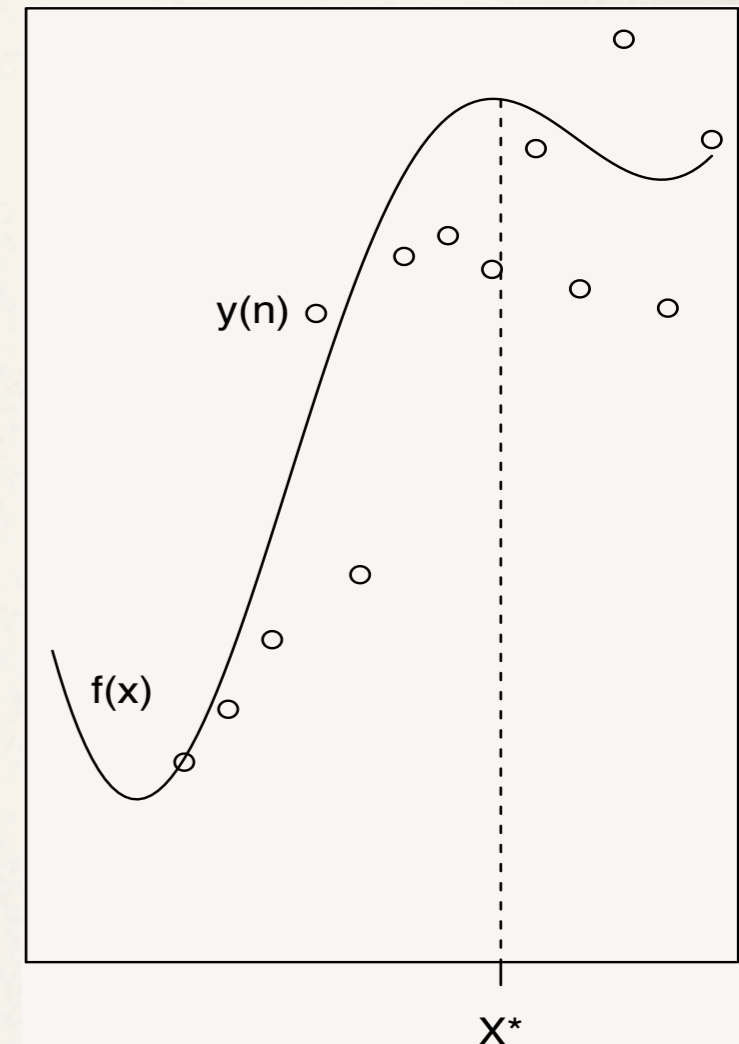
x



Algorithm

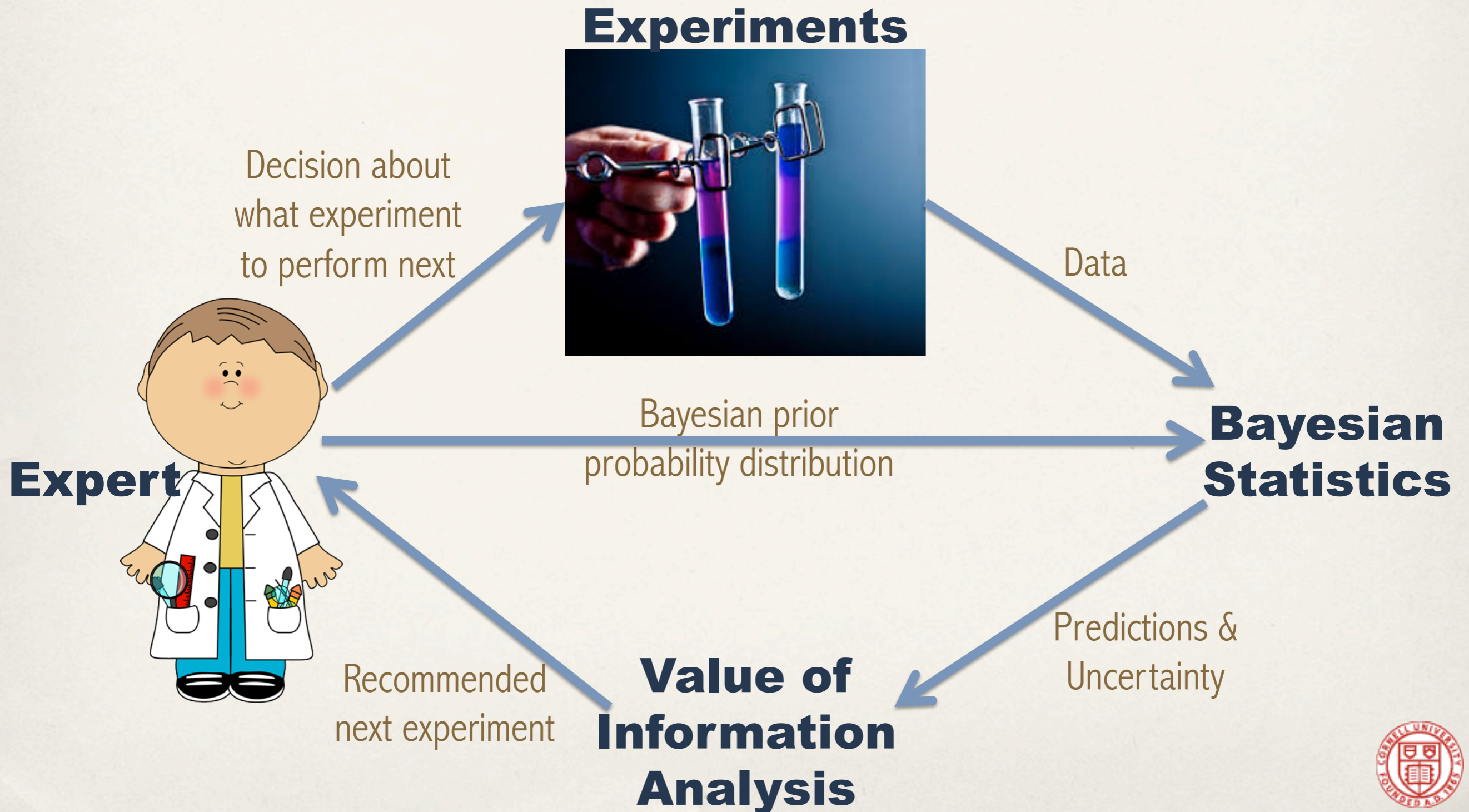


Objective function:
 $f(x)$



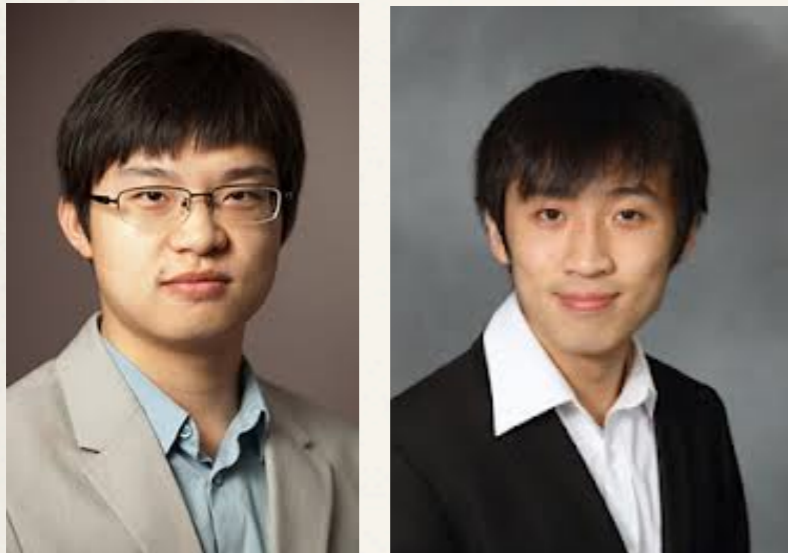
Objective function

We use Bayesian optimization to do trial & error in science



We use Bayesian optimization to make reversible fluorescent protein labels

- ❖ The Cornell team: Jialei Wang, Pu Yang.

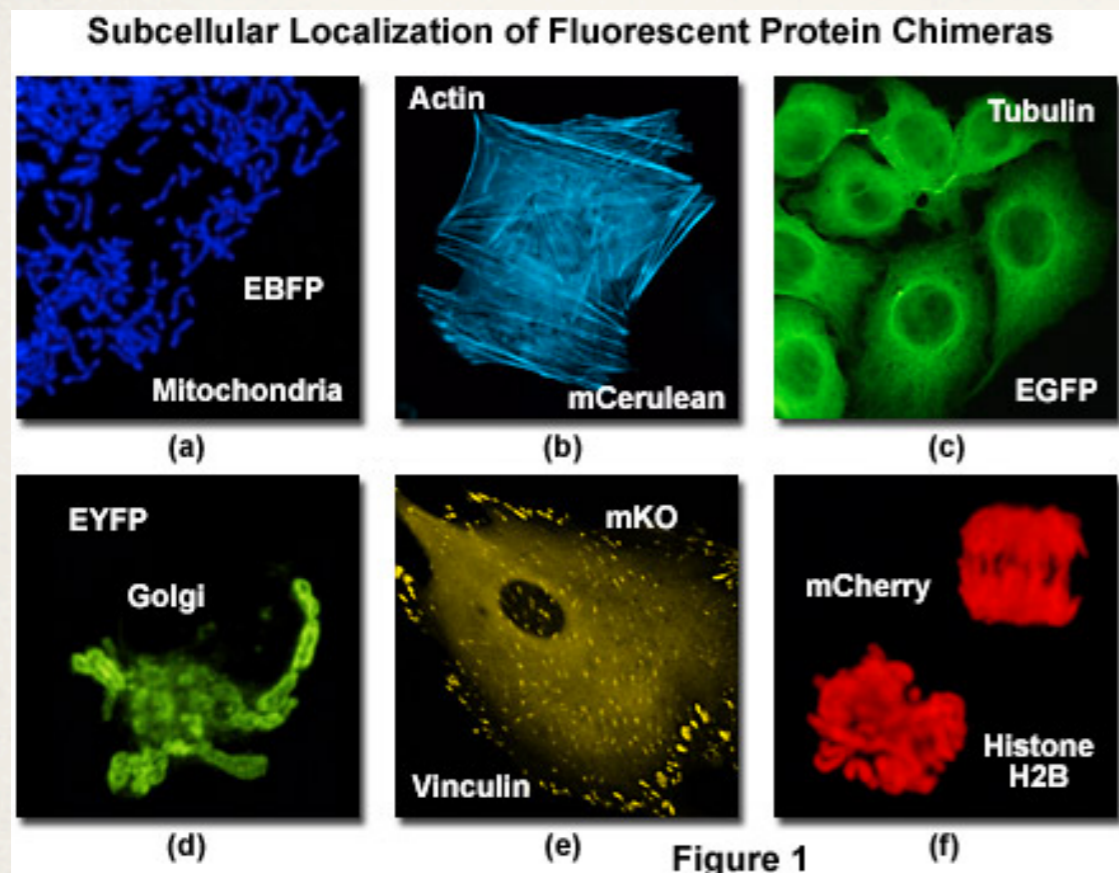


- ❖ Our scientific collaborators (all at UCSD): Mike Burkart, Nathan Gianneschi, Mike Gilson, Nick Kosa, Lori Tallorin.



Background: fluorescent protein labels are useful

- ❖ Scientists can attach fluorescent labels to proteins.
- ❖ This allows them to detect and track proteins as they move through living systems.



Source: www.olympusconfocal.com

- ❖ BUT, removing the labels is hard, or impossible.



Our goal is to make **reversible** fluorescent protein labels

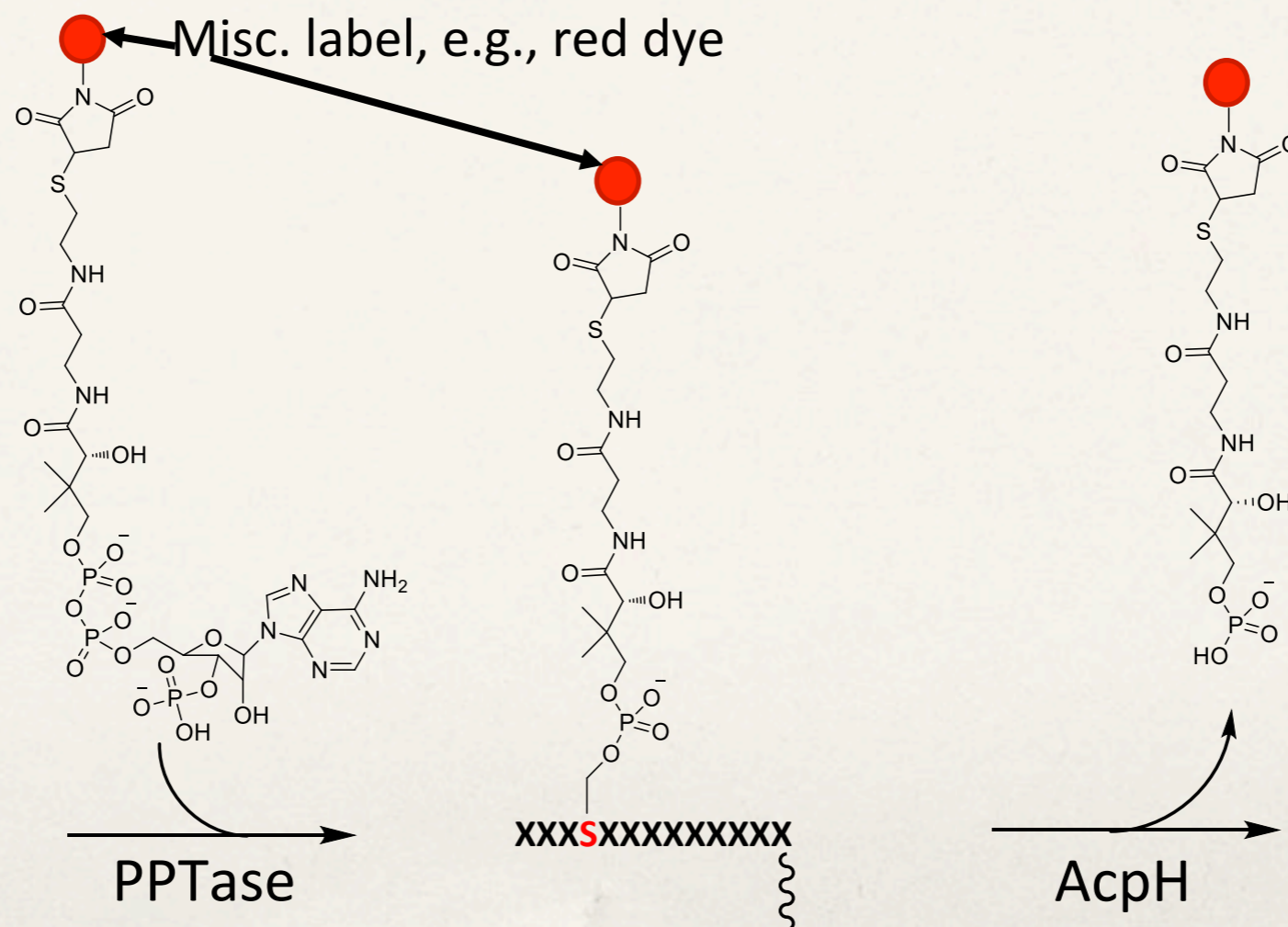
- ❖ If we could **remove** the fluorescent label from a protein, and then add it back, at will, it would be more useful.
- ❖ If we could attach not just fluorescent labels, but other kinds of labels that allowed us to manipulate them, e.g., by causing them to bind to beads, that would be even more useful still.



To make reversible protein labels, we need to find a peptide that allows two chemical reactions to happen

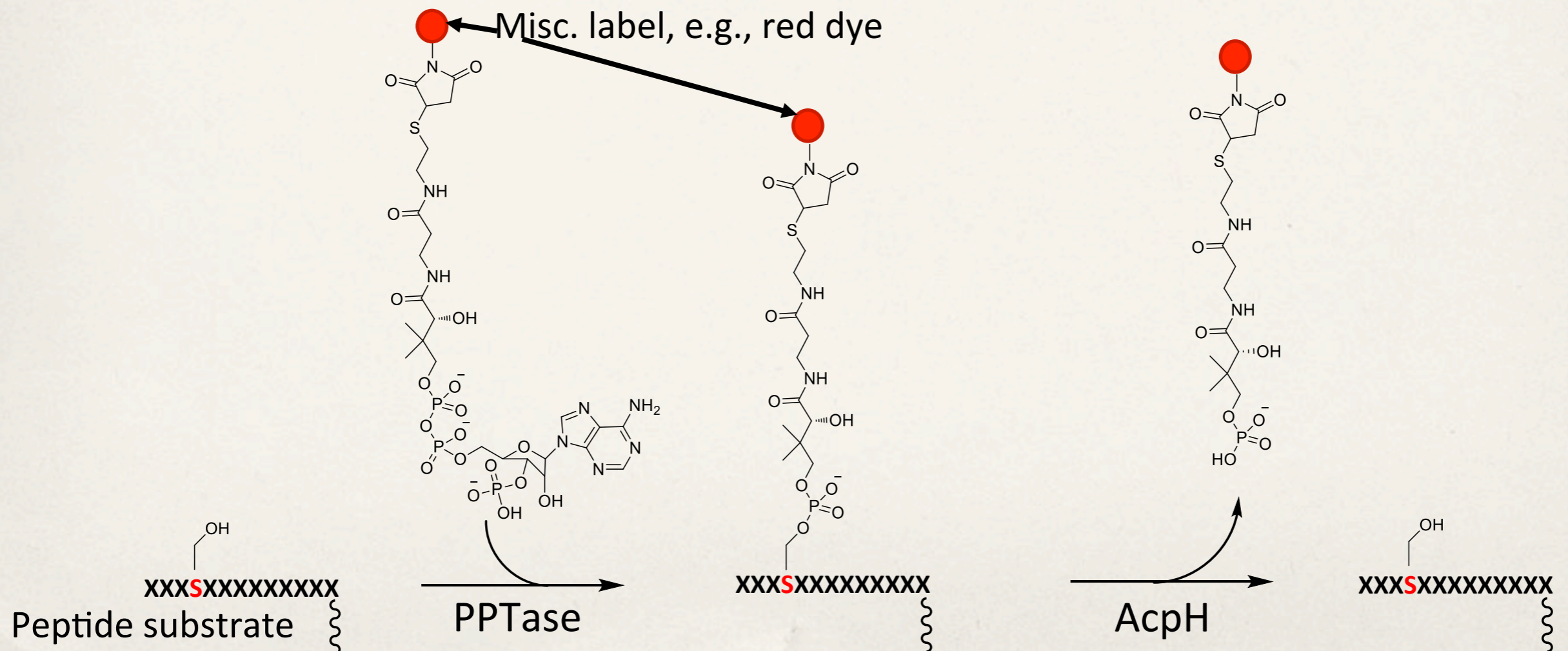
- ❖ Reaction 1 (with the enzyme PPTase) attaches the label to the peptide (at the “S”, short for “serene”).
- ❖ Reaction 2 (with the enzyme AcpH) removes the label.

- ❖ If we can find a peptide with this property, we’ll embed it within the protein to be reversibly labeled.



To make a reversible labeling system, we need to find a short hit

- ❖ If a peptide allows both chemical reactions to occur, we say it is a “hit”.
- ❖ For the reversible labeling system to be useful, the peptide needs to be short.
- ❖ Otherwise it will change the behavior of the protein in which it is embedded.



It is hard to find short hits; Math makes it easier.

- ❖ If a peptide allows both chemical reactions to occur, we say it is a “hit”.
- ❖ Hits are rare: about 1 in 10^5 among shorter peptides.
- ❖ Testing peptides is expensive & time-consuming: it requires reserving time on an expensive capacity-limited time machine, about 1 week’s worth of work by an experimentalist; and material costs.
- ❖ We test 500 peptides at time. 500 is much smaller than 10^5 .
- ❖ To help us, we have some known hits, obtained from natural organisms. They are too long to be used directly.

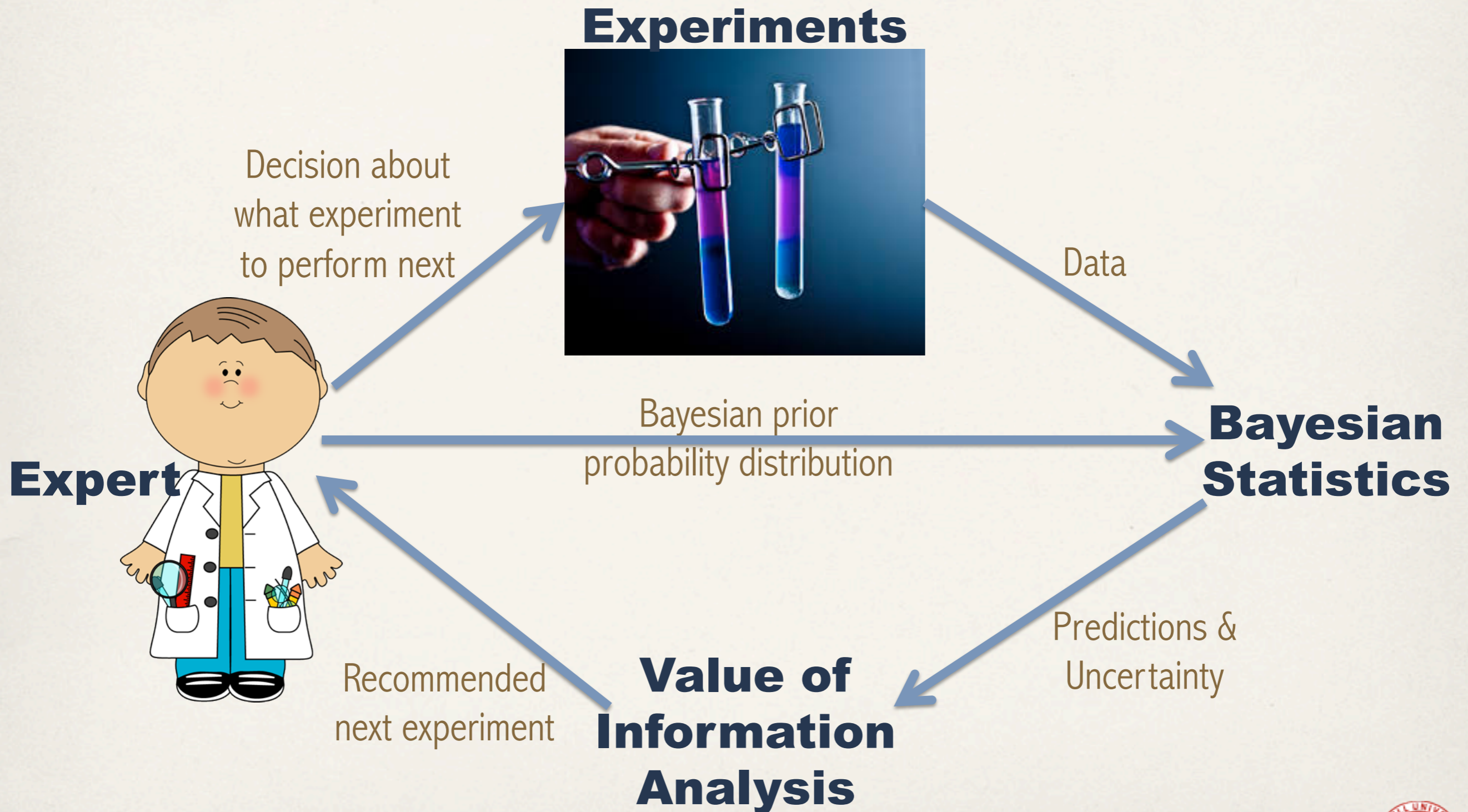


We reduce the experimental effort required to find minimal substrates

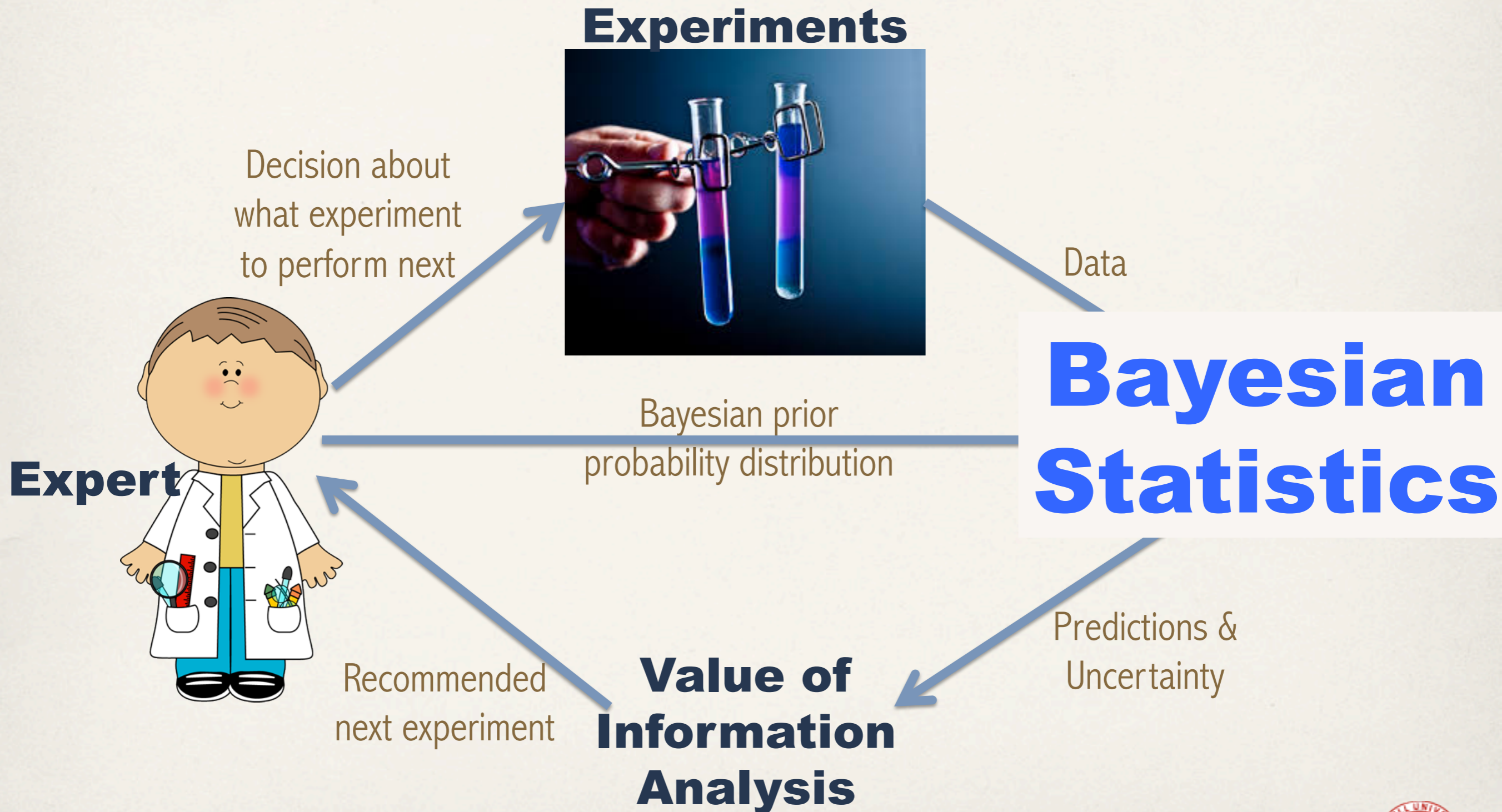
- ❖ We provide a method for Peptide Optimization with Optimal Learning (**POOL**).
- ❖ Our method has two parts:
 - ❖ Predict which peptides are “hits”, using Bayesian statistics.
 - ❖ Based on these predictions, recommend which peptides to test next, using value of information analysis.



Peptide Optimization with Optimal Learning (POOL)



First, we consider prediction.



We use Naive Bayes

- ❖ Naive Bayes is a statistical model often used for text classification (e.g., spam filters).
 - ❖ It is called “naive” because it makes a key independence assumption.
 - ❖ Although it is naive, it often works really well.
- ❖ We apply a variant of Naive Bayes to our problem, which is customized to include amino acids’ **location** within the peptide.



We use Naive Bayes

- ❖ We assume that reality is characterized by a pair of latent matrices, called $\theta^{(\text{hit})}$ and $\theta^{(\text{miss})}$, where columns of each matrix correspond to different positions within the peptide, and rows correspond to different types of amino acids.
- ❖ These latent matrices are unknown, but can be estimated from data.

- ❖ We further suppose that, for a peptide x ,

$$P(y(x) = 1 | x, \theta^{\text{hit}}, \theta^{\text{miss}}) = \frac{P(\text{hit}) \prod_i \theta_{i,x_i}^{(\text{hit})}}{P(\text{hit}) \prod_i \theta_{i,x_i}^{(\text{hit})} + P(\text{miss}) \prod_i \theta_{i,x_i}^{(\text{miss})}}$$

- ❖ Here, x is a peptide, x_i is the type of the amino acid at position i , $y(x)$ indicates whether x is a hit (1) or not (0), and $P(\text{hit})$ and $P(\text{miss})$ are prior estimates of the fraction of hits and misses in the population.

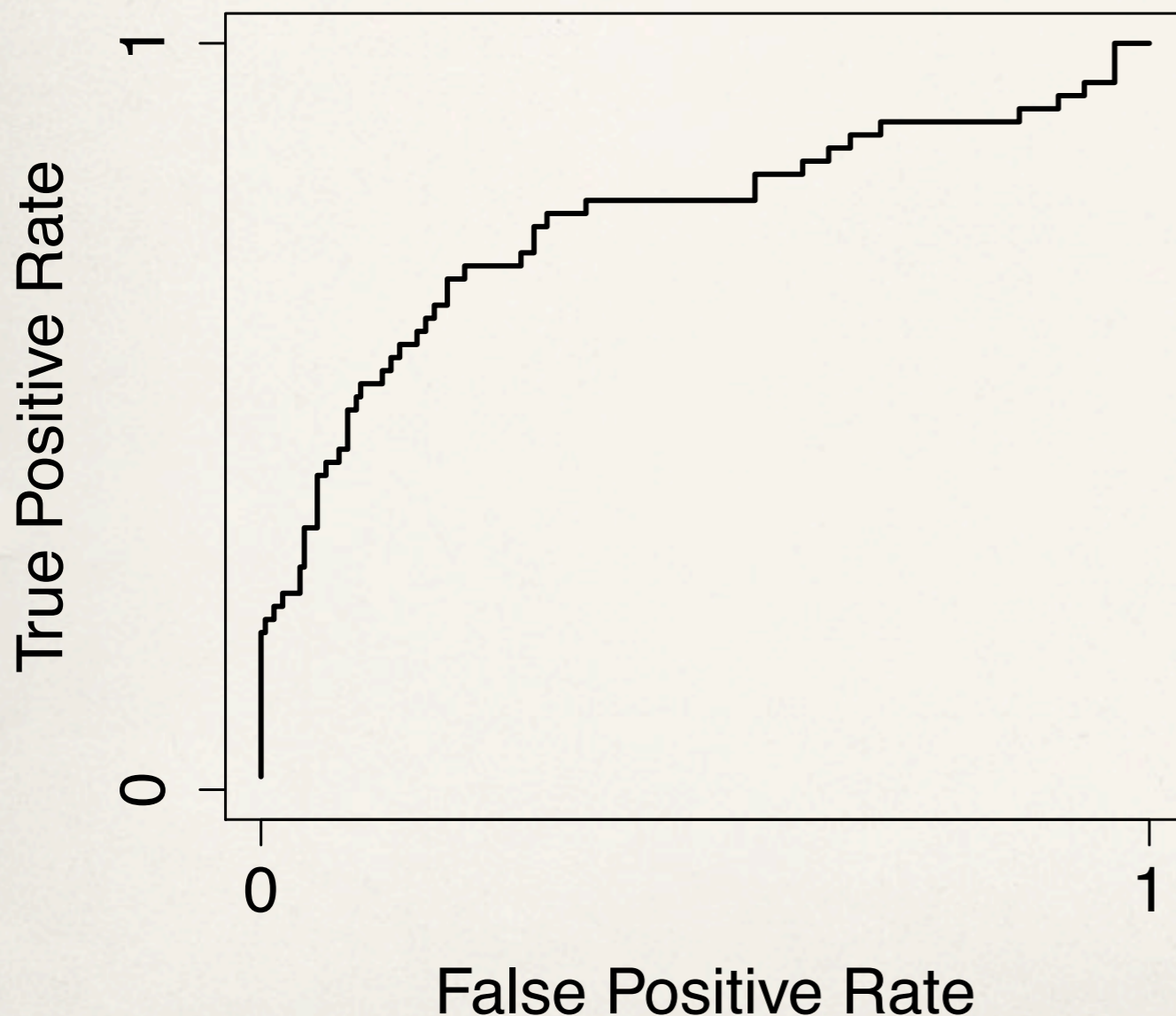


We use Bayesian Naive Bayes

- ❖ We put independent Dirichlet prior distributions on each column of the latent matrices $\theta^{(\text{hit})}$ and $\theta^{(\text{miss})}$.
- ❖ Our choices for the parameters of this prior are based on a biological understanding of the problem, discussions with our collaborators, and cross validation.
- ❖ Given training data $x^1, \dots, x^n, y(x^1), \dots, y(x^n)$, the posterior on the θ 's is also Dirichlet, and independent across i and j .
- ❖ To estimate the posterior probability of a hit, we can sample the thetas from the posterior, or calculate a single MAP estimate. The MAP estimate ignores uncertainty, but can be computed analytically.



This ROC curve suggests Naive Bayes performs reasonably well



- ❖ We have training data for approximately 300 peptides (most are misses.)
- ❖ True positive rate = % of hits labeled as hits.
- ❖ False positive rate = % of misses labeled as hits.
- ❖ Rates were estimated via leave-one-out cross-validation.



Now, we consider the choice of experiment

Experiments



Decision about what experiment to perform next

Data

Bayesian prior probability distribution

Bayesian Statistics

Expert



Recommendation for next experiment

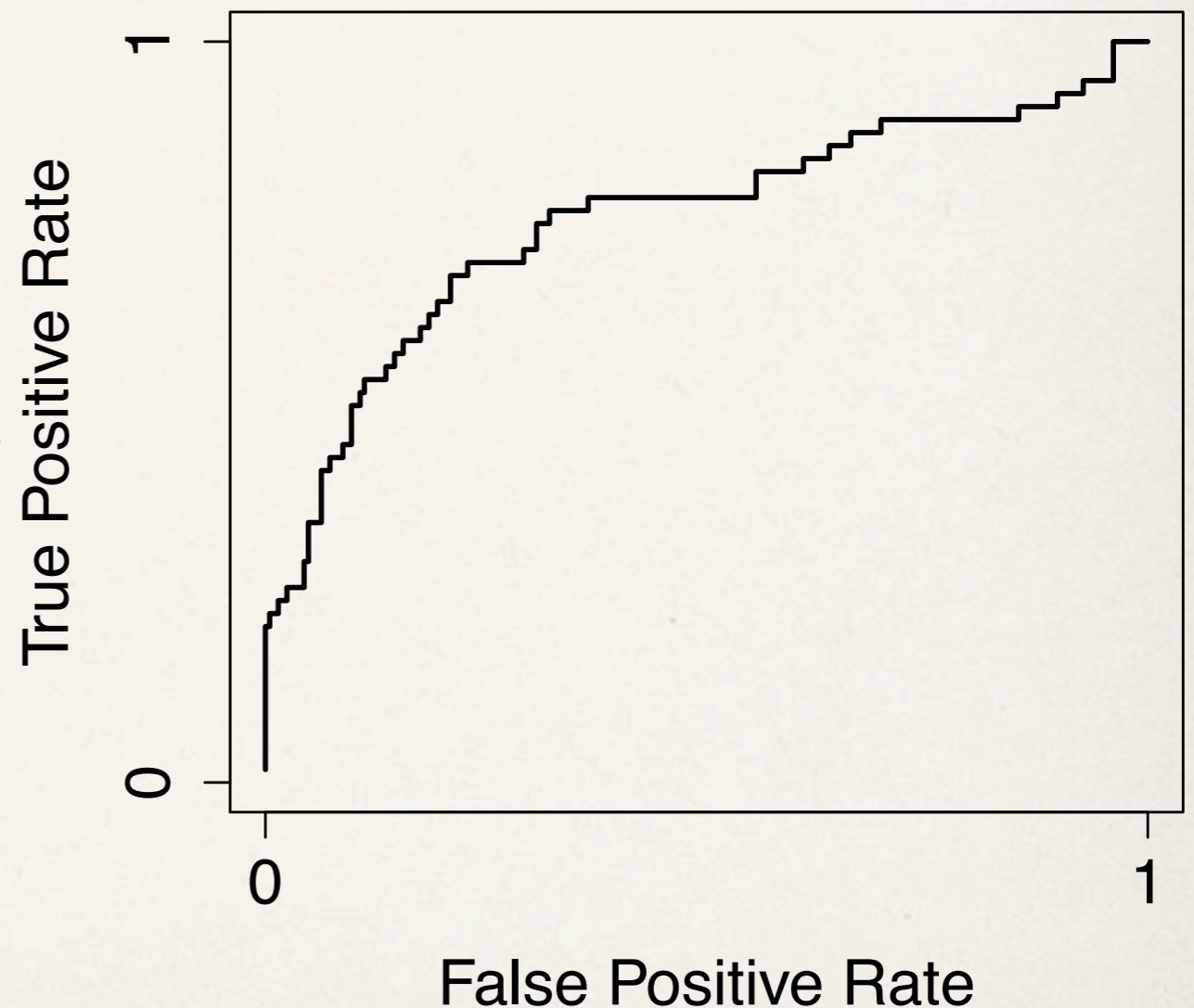
predictions & uncertainty

Value of Information Analysis



Given imperfect predictions, what should we test next?

- ❖ If predictions were perfect, we could just test the shortest peptide predicted to be a hit.
- ❖ Our predictions are not perfect.
- ❖ How should we decide what to test next?

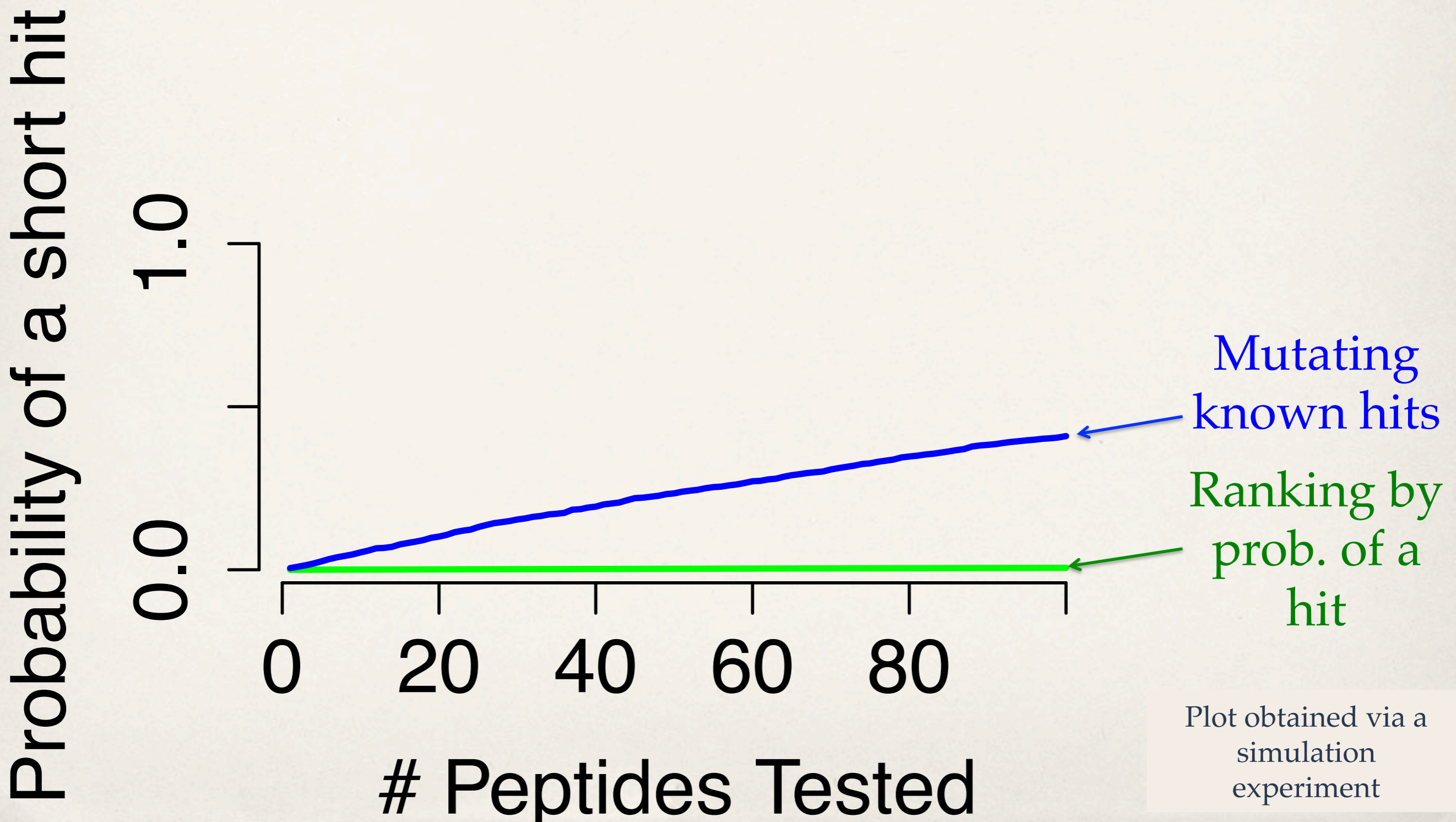


Ranking by probability of a hit does not work well

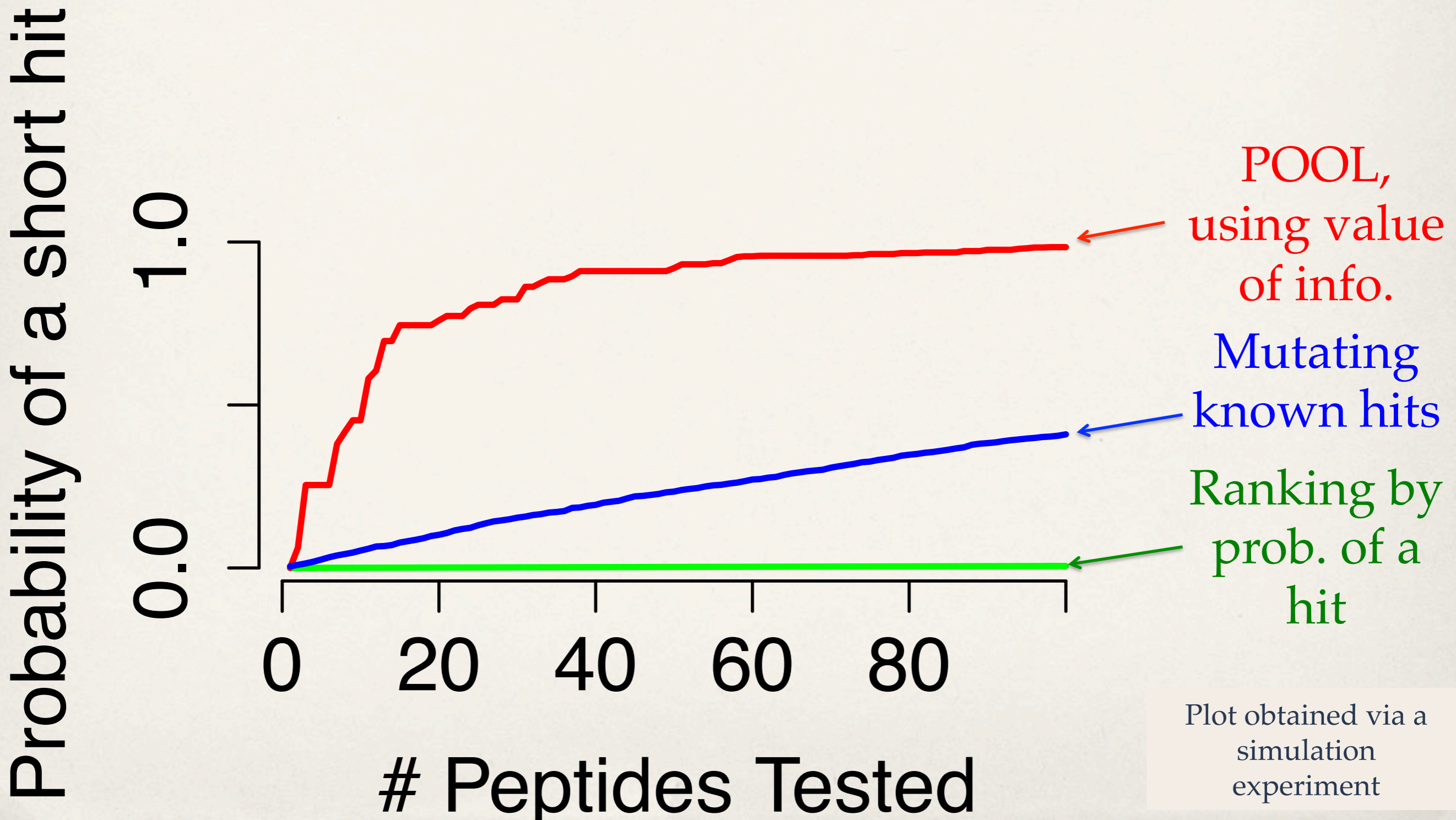
- ❖ One simple strategy is:
 - ❖ Select those peptides with length < 12 .
 - ❖ Rank them by predicted probability of a hit
 - ❖ Test the top 300.
- ❖ The tested peptides are very similar. If the first tested peptide is not a hit, the other ones probably aren't either.



Ranking by probability of a hit does not work well



Using value of information (VOI) works better



Value of Info. chooses the experiment that maximizes the probability we reach our goal

- ❖ Our goal is to find short hits.
- ❖ More specifically, our goal is:
 - ❖ Find at least one hit shorter than a target length b .
- ❖ We should run an experiment that maximizes the probability of reaching this goal.



The best experiment is the solution to a combinatorial optimization problem

- ❖ This can be formulated as this combinatorial optimization problem:

$$\max_{S \subseteq E: |S| \leq k} P(\text{at least one short hit in } S)$$

- ❖ Notation:
 - ❖ E is the set of all peptides.
 - ❖ S is the set of peptides to test.
 - ❖ k is the number of peptides we can test in one experiment. Typically, k is between 200 and 500.
 - ❖ A “short hit” is a hit whose length is less than b.



We can't solve this exactly, so we approximate the solution using a greedy algorithm

- ❖ This combinatorial optimization problem is very challenging : The size of the set $\{S \subseteq E : |S| \leq k\}$ is $|E|$ choose k . If $b=15$ and $k=500$, this is 10^{19} choose 500.
- ❖ Instead, we build up the set S of peptides to test in stages.
- ❖ In each stage, we find the single peptide to add that maximizes the probability of reaching our goal:

$$\max_{e \in E \setminus S} P(\text{at least one short hit in } S \cup \{e\})$$

- ❖ We then add e to S and repeat, until S has $k=500$ peptides.



The greedy algorithm has an approximation guarantee

Lemma: $P^*(S)$ is a monotone submodular functions of S .

Proposition: Let $\text{OPT} = \max_{S \subseteq E: |S| \leq k} P^*(S)$, and let GREEDY be the value of the solution obtained by the greedy algorithm. Then

$$\frac{\text{OPT} - \text{GREEDY}}{\text{OPT}} \leq 1 - 1/e$$

- ❖ In the above, $P^*(S) = P(\text{at least one short hit in } S)$.
- ❖ The proposition follows from the lemma & a result from Nemhauser, Wolsey, Fisher '78.
- ❖ This result is similar in spirit to results obtained in Y. Chen & A. Krause, "Near-optimal Batch Mode Active Learning and Adaptive Submodular Optimization," ICML 2013.



We can implement the greedy algorithm efficiently

- ❖ The greedy optimization step can be shown to be equivalent to

$$\arg \max_{e \in E \setminus S} P(y(e) = 1 | y(x) = 0 \forall x \in S)$$

- ❖ We can compute this probability by treating all peptides in S as misses, and re-training our model. If we then use a MAP estimate, this probability decomposes over the amino acids, and can be optimized efficiently.



Here is the intuition why this approach works better than “rank by prob. hit”

- ❖ Finding the the single peptide to add that maximizes the probability of reaching our goal:

$$\max_{e \in E \setminus S} P(\text{at least one short hit in } S \cup \{e\})$$

- ❖ Is equivalent to:

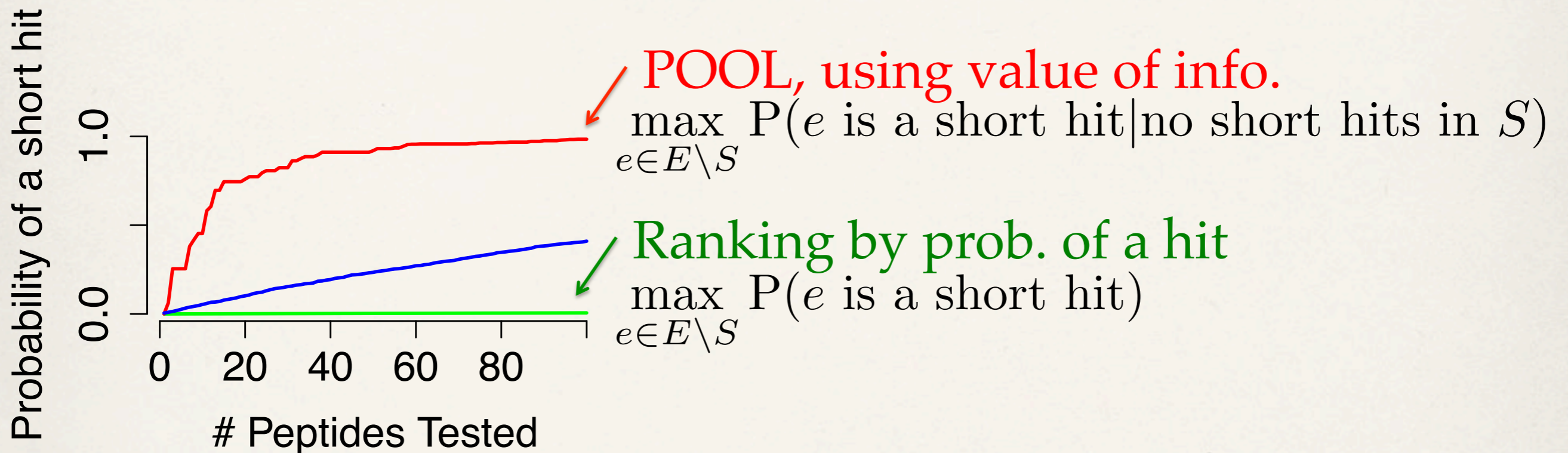
$$\max_{e \in E \setminus S} P(e \text{ is a short hit} | \text{no short hits in } S)$$

- ❖ Compare this to the “rank by prob. hit” approach

$$\max_{e \in E \setminus S} P(e \text{ is a short hit})$$



VOI works better because its peptides are more diverse

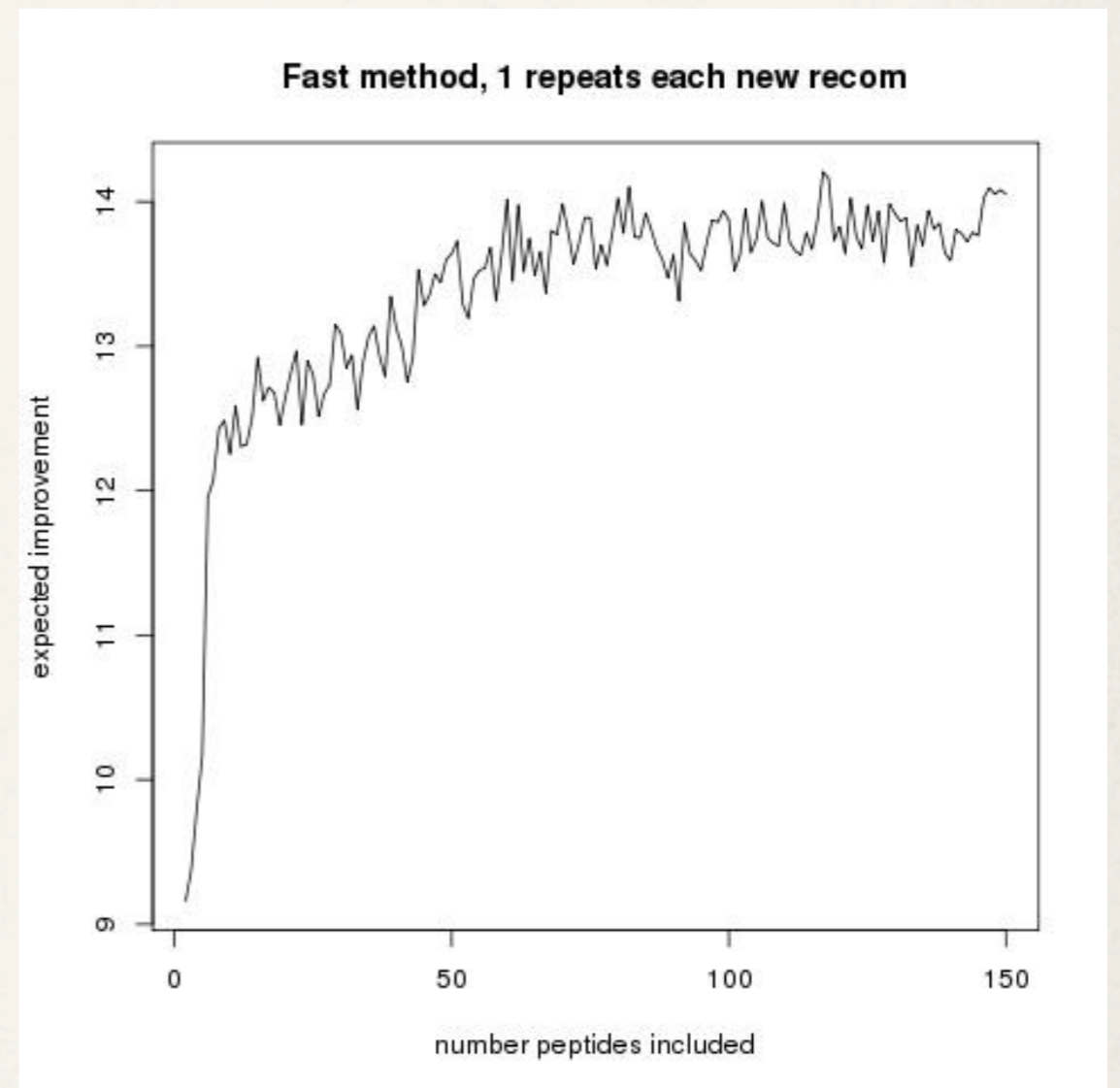


- ❖ Peptides added using the value of information approach tend to be **different** from those already in S .
- ❖ Its recommendations are more **diverse**.



We are greedily optimizing $P(\geq 1 \text{ short hit})$ with one tweak to make real recommendations

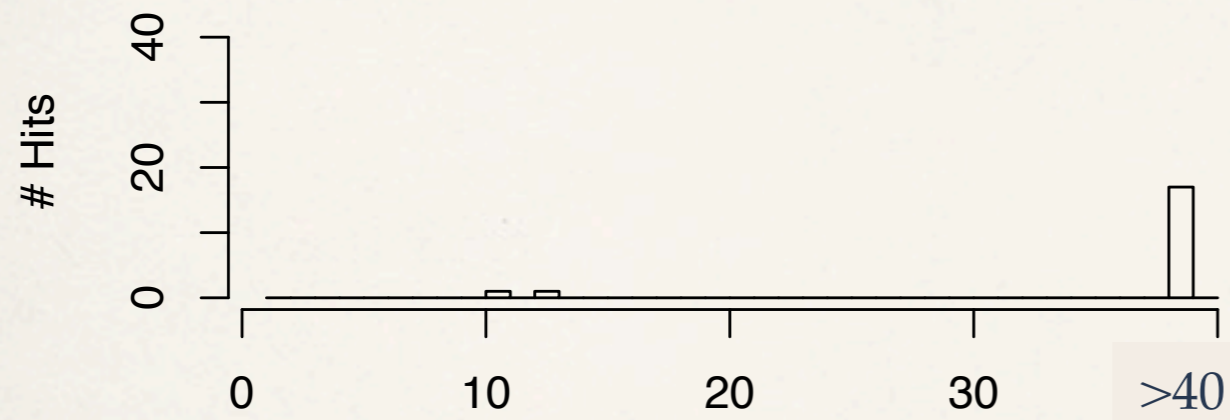
- ❖ We have used the following approach in recommending experiments to our collaborators.
- ❖ We pre-select a random sequence of lengths a^1, \dots, a^k strictly less than b , and require that the n^{th} peptide selected has length less than a^n .
- ❖ We then apply the greedy probability of improvement algorithm.
- ❖ This makes the hits shorter, without hurting $P(\geq 1 \text{ short hit})$.



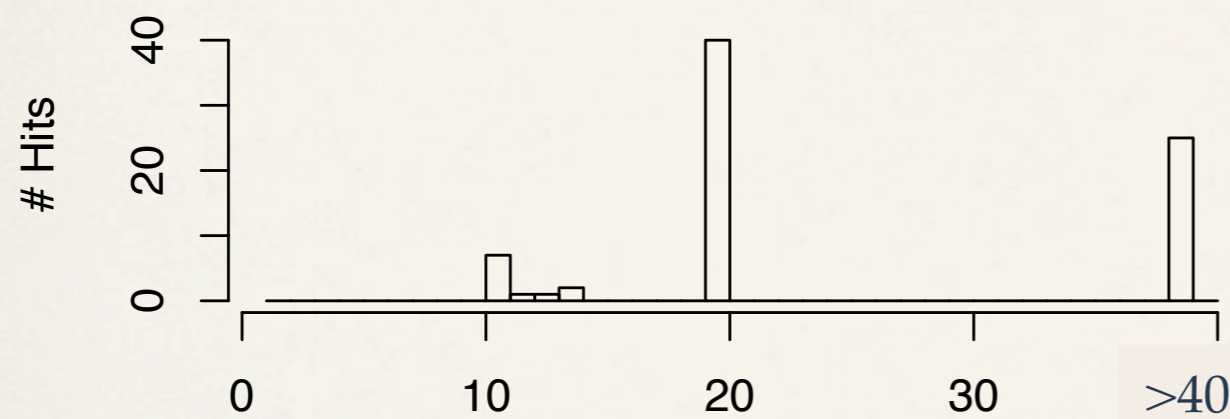
Expected improvement as a function of $|S|$, estimated via Monte Carlo.



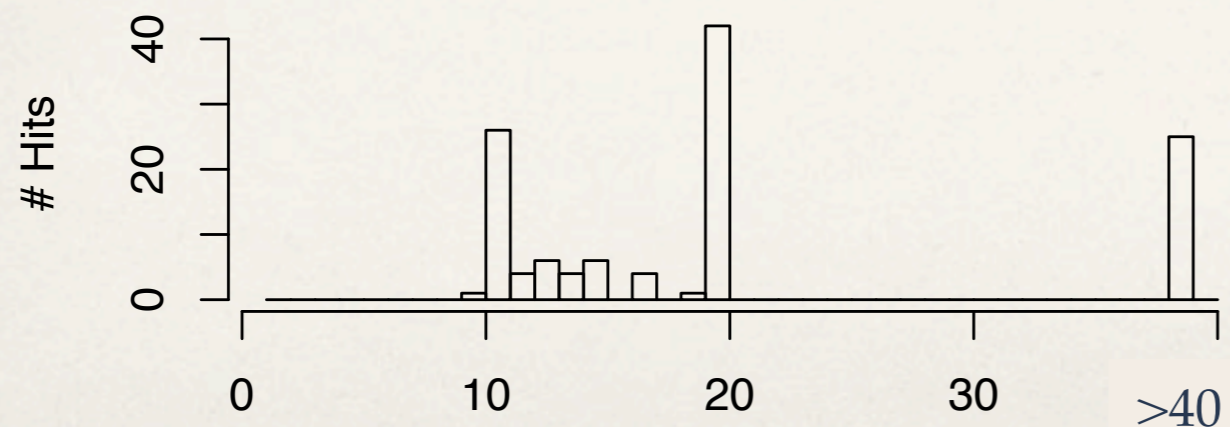
We have found novel short peptides using this method



Training Set
Length of shortest hit: 11



After 1 round of POOL
Length of shortest hit: 11



After 2 rounds of POOL
Length of shortest hit: 10

Peptide Length



Conclusion

- ❖ We have developed an optimal learning method for finding minimal peptide substrates.
- ❖ This method has found hits shorter than the shortest previously known.
- ❖ This approach can be applied to other kinds of problems.
- ❖ This approach aims to:
 1. Reduce the experimental effort required to reach a goal.
 2. Increase the chance of achieving a goal within a given experimental budget.



Additional slides



We can use ideas from Bayesian optimization to do trial & error in science

- ❖ Bayesian optimization goes back to [Kushner, 1964], with important developments due to Mockus in the 70s and 80s. It is a kind of Bayesian experimental design.
- ❖ In Bayesian optimization, we
 - ❖ 1. Build a Bayesian statistical model of the objective function based on previously collected data.
 - ❖ 2. Suggest experiments to run next with a large value of information [Howard 1966].
- ❖ Bayesian optimization can be seen as a surrogate optimization technique. Other kinds of surrogate optimization should also be useful for trial & error.
- ❖ Bayesian optimization is in a larger class of methods I call “optimal learning” methods.



Using VOI to optimize $P(\geq 1 \text{ short hit})$ has a shortcoming

- ❖ Under our Naïve Bayes model, it is usually possible to increase $P(\text{hit})$ by increasing the peptide's length.
- ❖ Thus, the experiments that maximize $P(\geq 1 \text{ short hit})$ tend to have length $b-1$.
- ❖ However, a hit strictly shorter than $b-1$ would be even better.
- ❖ To allow us to find such strictly shorter peptides, we might consider an alternate goal: expected improvement.



Optimizing expected improvement would fix this

- ❖ Let $f(x)$ be the length of peptide x .
- ❖ $f^*(S) = \min_{x \in S: y(x)=1} f(x)$ is the length of the shortest hit found.
- ❖ Define the **expected improvement** for testing S as:
$$EI(S) = E[(b - f^*(S))^+]$$
- ❖ An S that maximizes $EI(S)$ could contain peptides shorter than $b-1$.



Efficiently optimizing expected improvement is ongoing work

- ❖ Solving $\max_{S \subseteq E: |S| \leq k} \text{EI}(S)$ exactly is very challenging.
- ❖ $\text{EI}(S)$ is also a monotone submodular function, and so the greedy algorithm also has an approximation guarantee.
- ❖ However, actually finding the single peptide to add that maximizes the expected improvement is itself extremely difficult.
- ❖ We are currently using an integer program to do this, but results are pending.

