

Experimental (T)error: Dealing with data in an uncertain world



Terry R. Stouch
Science for Solutions, LLC
Consulting in Drug discovery and design

John D. Chodera
Computational Biology Program, Memorial Sloan-Kettering Cancer Center
<http://www.choderalab.org>

Motivation

How can we, as modelers, deal with data error in a useful and reasonable manner?

How should error affect our interpretation and use of the data?

Our premise:

Realistic assessments of error can save you **time, money, and sanity by ensuring you are chasing signal, rather than fitting the noise**

The unfortunate truth:

Our field currently lacks standard practices for dealing with error.

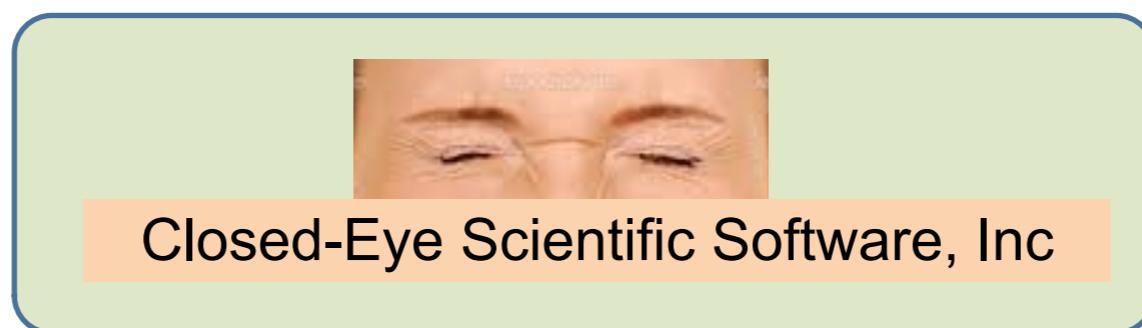
Introduction to error

- Endemic
- Ignored
- Hidden
- Substantial
 - Could obviate the use of the data
- High precision / low error is expensive
 - Tremendous care and expertise might be required

Introduction to error

- Sometimes people – users of the data - just don't care!
 - “I don't care if the measurements aren't quantitative, I need to see several significant figures”
 - “I don't care if the answer is right, as long as it's fast”
 - ~”If someone tells me they don't trust a number, I'll tell them to get out of my office and not return until they get it right.”
 - ~”Just give me the numbers, I'll know how to interpret the data.”
 - ~”Give me your data, I can get a model that will fit it to 100%”

“It's all automated now. The QSAR problem has been solved”



The up side of this talk!

- Understanding error helps us avoid and accommodate it resulting in
 - Better decisions
 - Better models
 - More appropriate use of the data and the models
- It might lead to better experiments
- Models can replace high error experiments
- Experimentalist could concentrate on producing high quality and novel data



Outline

What are the sources of experimental error?

How do we measure and report error?

How can we estimate error if not measured or reported?

When is it safe to use literature data?

What level of error is needed for data to be useful?

How do we propagate error?

Can we really get rid of outliers and feel good about ourselves in the morning?

When is it better to stick with a model than to collect data?

Some rules of thumb (for when you have no other information)

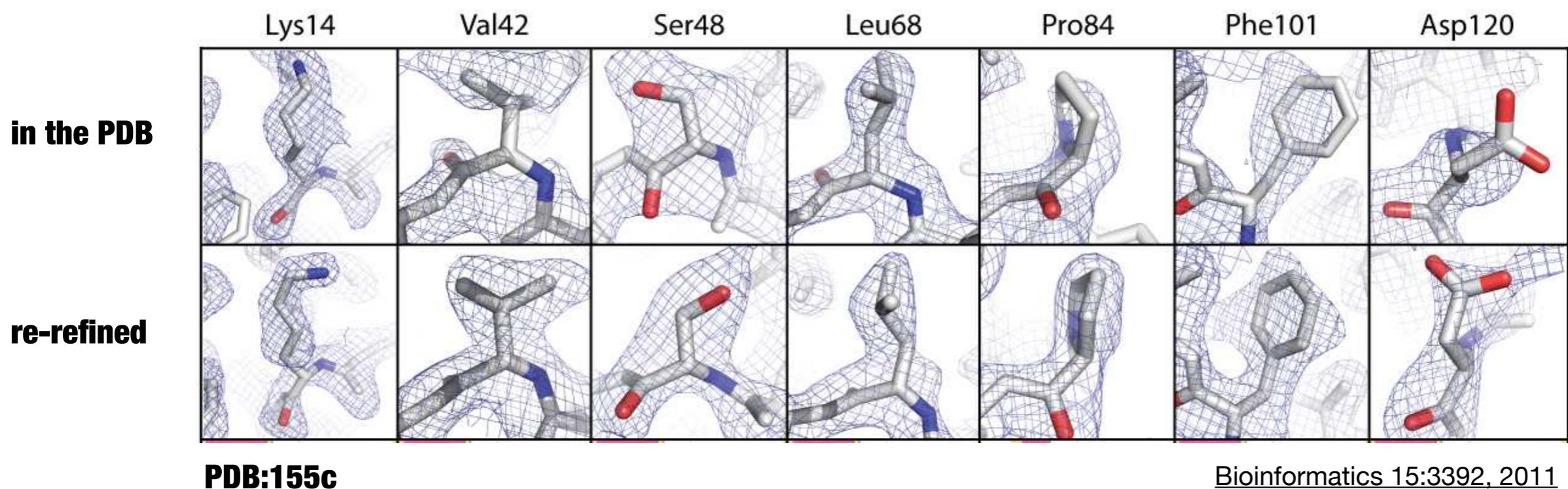
Pitfalls, catastrophes, and tales of woe

What we **aren't** going to cover

Errors in structural data (X-ray and NMR)

Some brief warnings:

- X-ray and NMR structures are just **models** derived from data
- often the models don't capture important aspects of the data
- all data-derived models have **uncertainty**: know how to quantify this
- refinement **errors** can be more frequent than you might expect



Genesis of the data

- Likely it was not generated for long-range data mining
- Largely generated for a specific need or program at a specific time
 - Consistency might be greater internal to a project and during a particular temporal period

Different sources of error to watch for

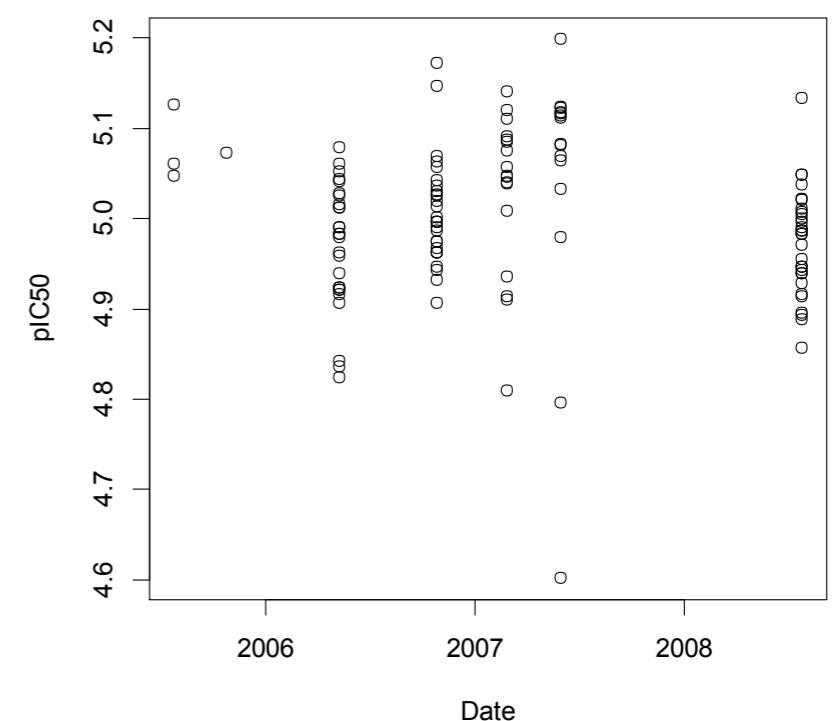
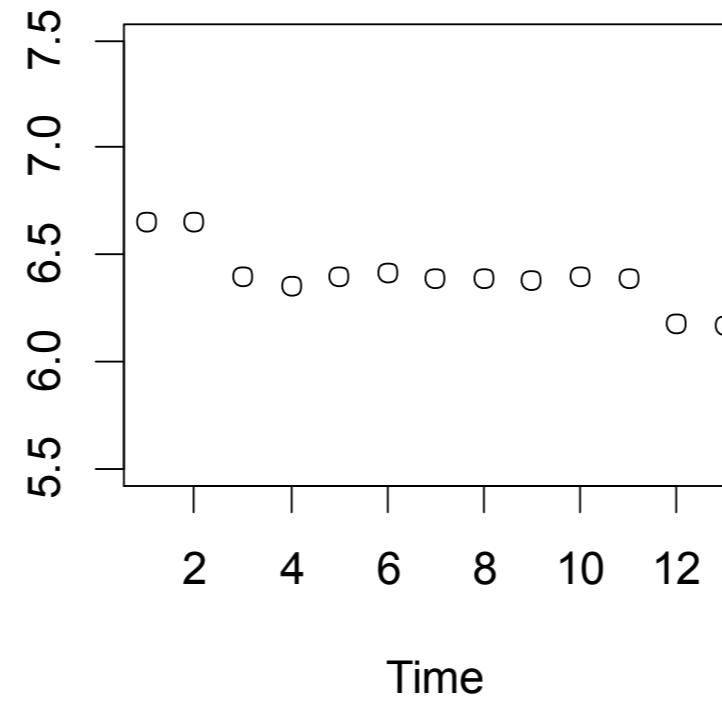
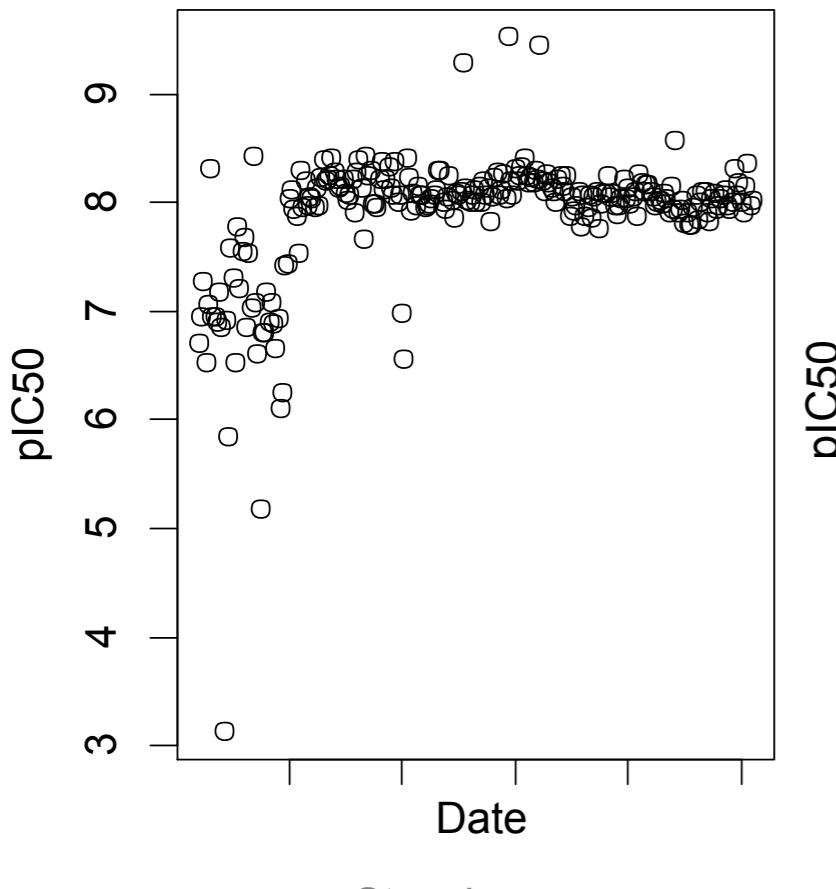
- Concentration of *desired* substance
 - Purity
 - Age
 - Stereochemical purity
 - Effects of impure fraction
 - Weighing precision
 - Solubility
 - Reliability of the determination
- Source of material
 - In house, recent synthesis vs. purchased compound collection
- Source of the data
 - HTS
 - “Program” data

Different sources of error to watch for

- Different Standard Operating Procedures (SOP) same company
 - The same ultimate endpoint arrived at with different methods
 - Eg. LC/UV vs LC/MS for measurement
 - Weight / volume to determine concentration vs. MS or NMR
- Same SOPS same company but different: labs, solutions, samples, days, equipment, technicians, water, humidity, temperature, lab coats, time of day,
 - Veterans of assays: expect 2X variation at least, 5X not unexpected, 10X no contention
- Different SOPs different institutions
 - Can it be any better than same SOP same company different labs?

Different sources of error to watch for

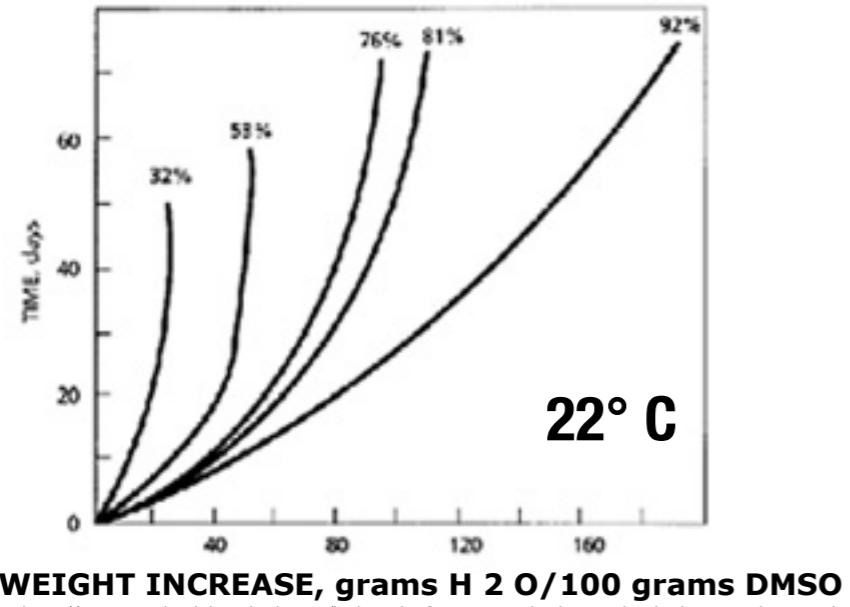
- How are standards used?
 - When bad values are obtained, is that plate / series ignored? Is the data still saved?
 - If the data is saved, is it flagged? Can data from ‘bad’ plates be identified?



An example source of unanticipated (t)error: DMSOh No!

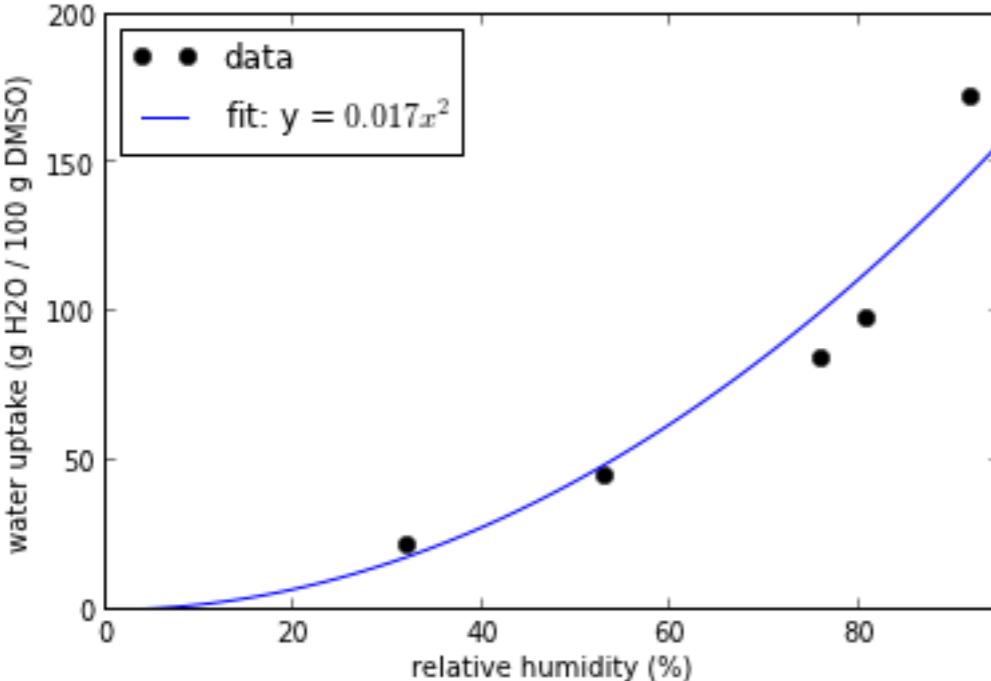
Most compound libraries are stored as stocks in 100% DMSO.

But DMSO is incredibly hygroscopic! Water uptake depends on relative humidity.

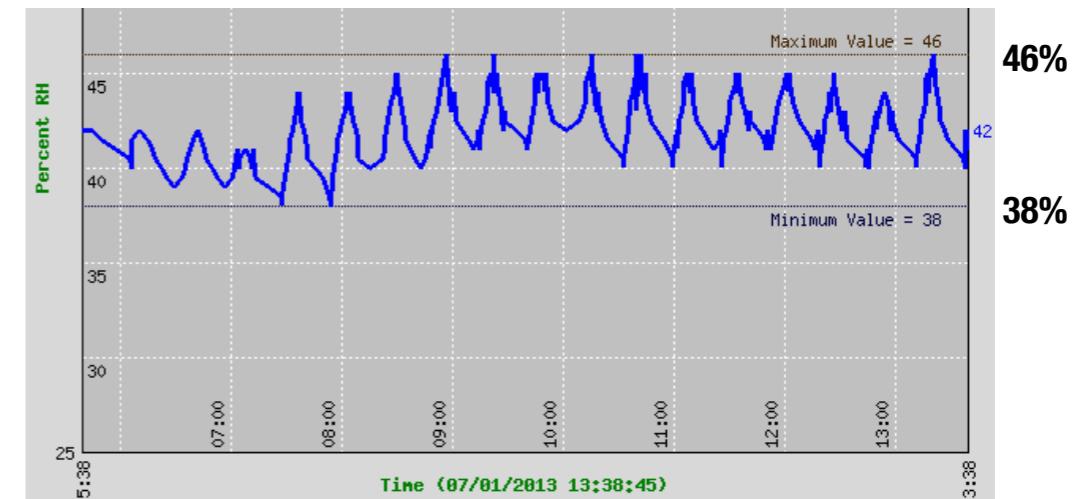


WEIGHT INCREASE, grams H₂O/100 grams DMSO

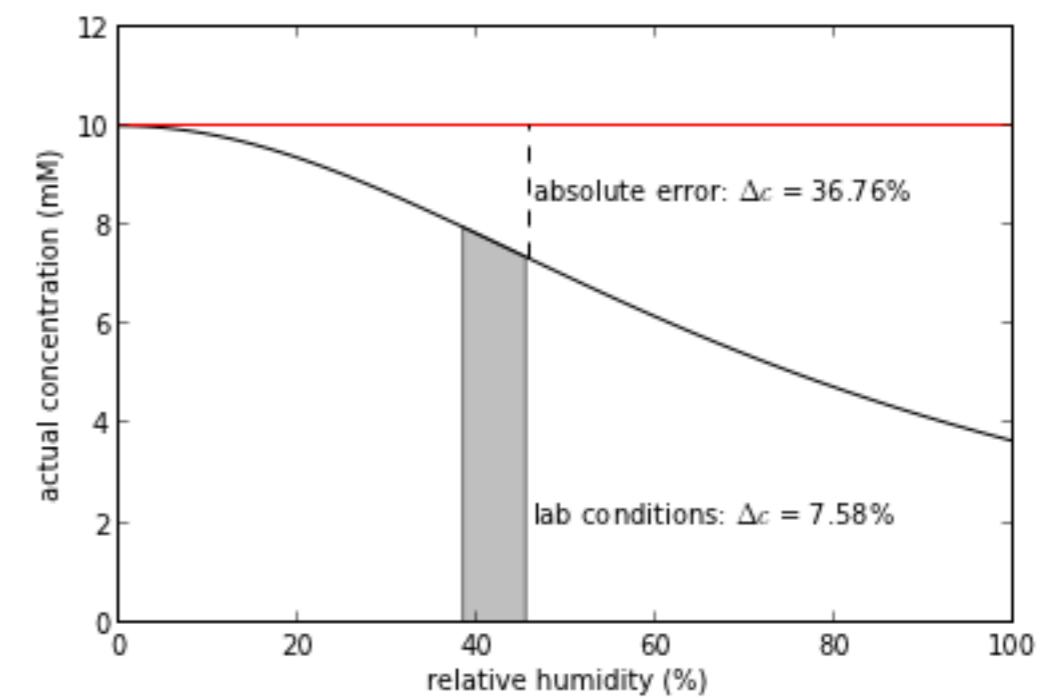
From <http://www.gaylordchemical.com/index.php?page=101b-dmso-physical-properties#101b-4>



If we left DMSO stock container open, absolute concentration error ~37% and daily fluctuations ~7.6%



relative humidity fluctuations in my ZRC 17th floor lab



Different sources of error to watch for

- Is the number *derived* (from unseen raw data)
 - What is the error in the derivation?
 - If IC₅₀, then are the diagnostics of the curve available?
- If concentration is involved in the endpoint, how is it determined?
 - Weight and volume? Weighing errors are a problem – typically 2X, often more
 - Mass spec?
 - NMR?
- Omission of metadata or odd metadata (could imply a problem)
 - Hill coefficients, Z-prime, confidence interval ratio, purity, dynamic range qualifiers
 - Odd dates, extreme age (1945 for date of synthesis)
- Inappropriate reporting of results
 - *Don't blame the data, blame the way it was presented*

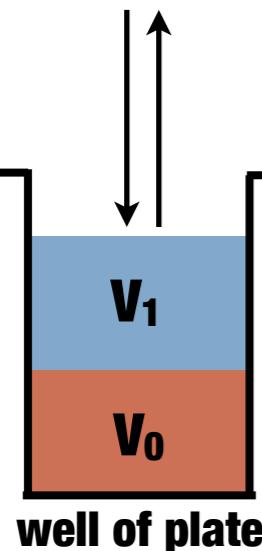


Science for Solutions, LLC

What is experimental error and where does it come from?

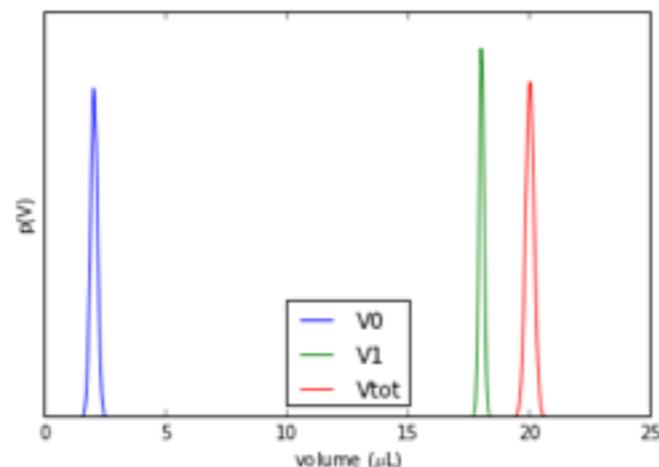
fluorescence measurement

assay mix
compound



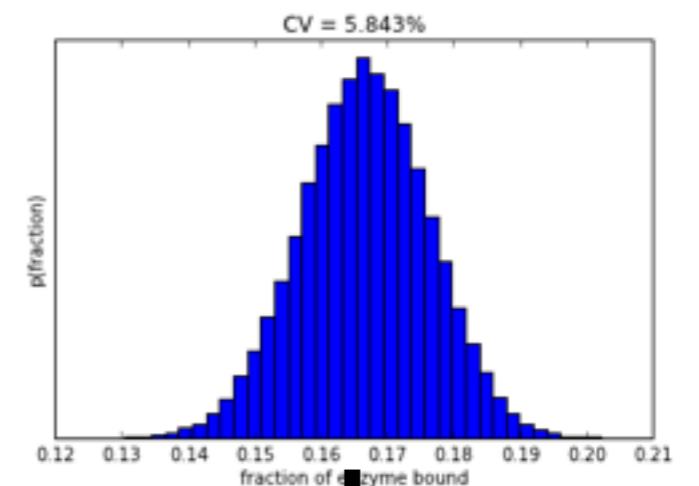
$$\begin{aligned}V_0 &= 2 \mu\text{L} \\V_1 &= 18 \mu\text{L} \\K_d &= 100 \text{ nM} \\C_0 &= 10 \text{ mM} \\C_1 &= 1 \mu\text{M}\end{aligned}$$

finite pipetting precision

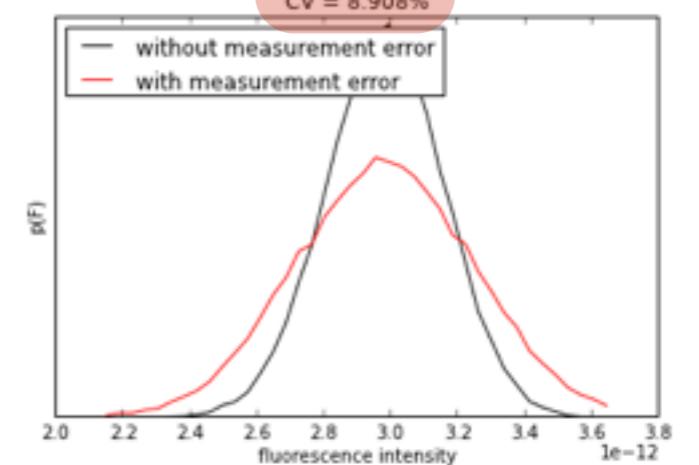


| BIOMEK FX ^P PIPETTING PERFORMANCE SPECIFICATIONS | | | | |
|---|------------------------------|-----------------------|--------------|---------------|
| SPAN-8 SYSTEMS | | | | |
| Transfer Volume | Span-8 Syringe Volume | Tip Types | Accuracy ± % | Precision < % |
| 0.5 μL | 250 μL | P20, Fixed 60 mm | 5 | 10 |
| 1 μL | 250, 500, 1000 μL | P20, P50, Fixed 60 mm | 3 | 7 |
| 5 μL | 250, 500, 1000 μL | P20, Fixed 60 mm | 3 | 5 |
| 10 μL | 500 μL | P50, P250 | 3 | 5 |

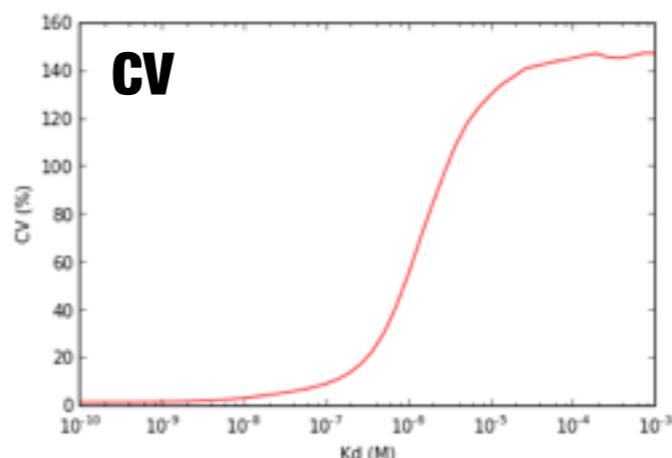
variation in protein:ligand complex concentration



CV = 8.908%



signal broadening due to measurement error



Is the expected CV small enough to be useful?

Modeling an experiment can ensure it will yield useful data or uncover unexpected issues

IPython tutorial available at https://github.com/choderalab/cadd-g_2013

The central limit theorem is why errors often look Gaussian

**Suppose X_n is a random variable from ANY distribution with finite variance.
The sample mean (average) of a few samples will always look Gaussian!**

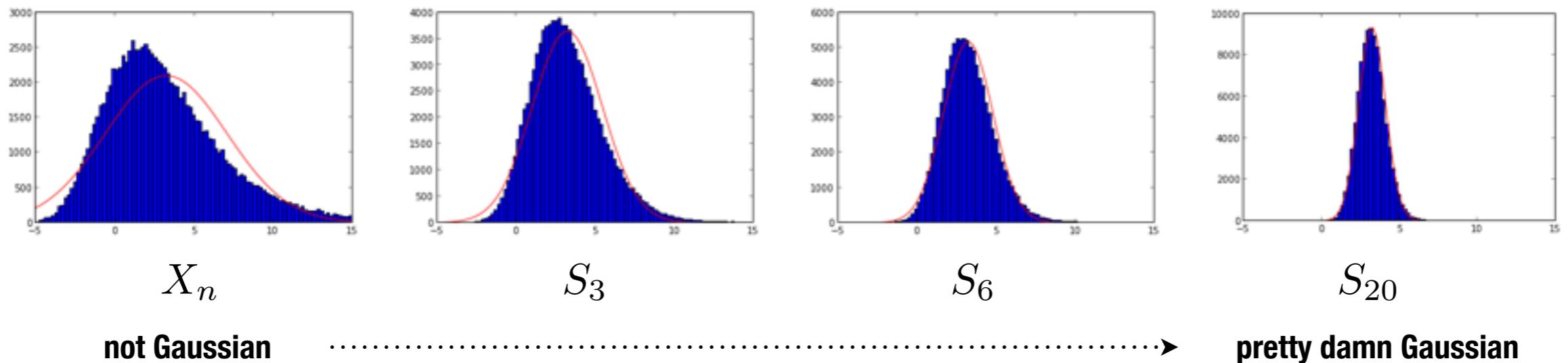
central limit theorem (CLT): $\frac{1}{N} \sum_{n=1}^N X_n \rightarrow N(\mu, \sigma^2/N)$

sample mean

$$E[X_n] = \mu$$

$$var[X_n] = \sigma^2$$

This happens with surprisingly few measurements, even for crazy distributions!



Many instruments operate like this (e.g. plate reader that uses 10 flashes of Xe lamp)

Reporting the error

Rules of error reporting:

- Report only **one significant digit** in the error estimate
- Report the measurement only to this decimal place
- Report what kind of error estimate you are reporting for numbers or error bars

Examples

YES: 4.1 ± 0.1 or $4.1(1)$

YES: 4.123 ± 0.002 or $4.123(2)$

NO: 4.123 ± 0.1 or 4.123 ± 0.123

YES: standard error of the mean

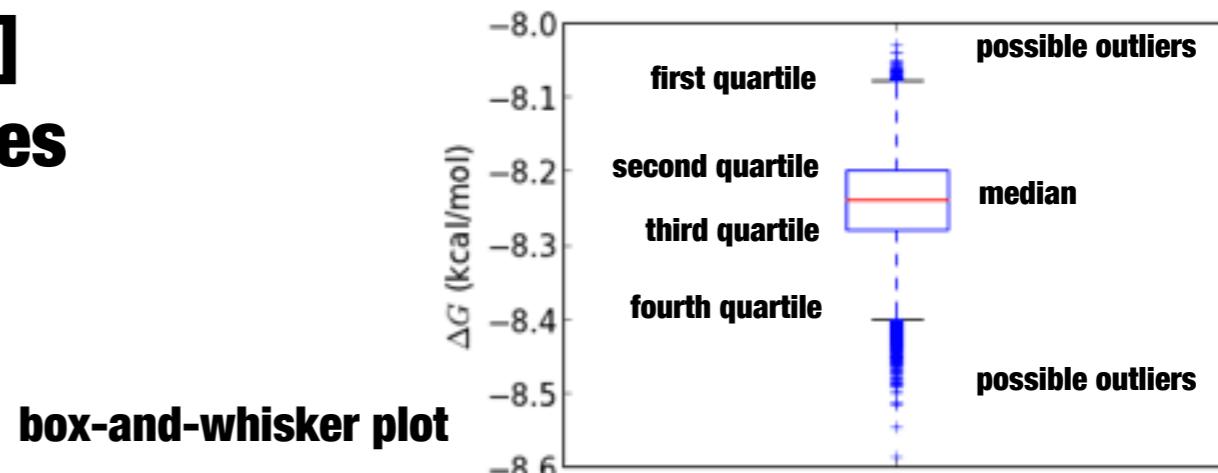
YES: 95% confidence interval

NO: meaning of error not specified

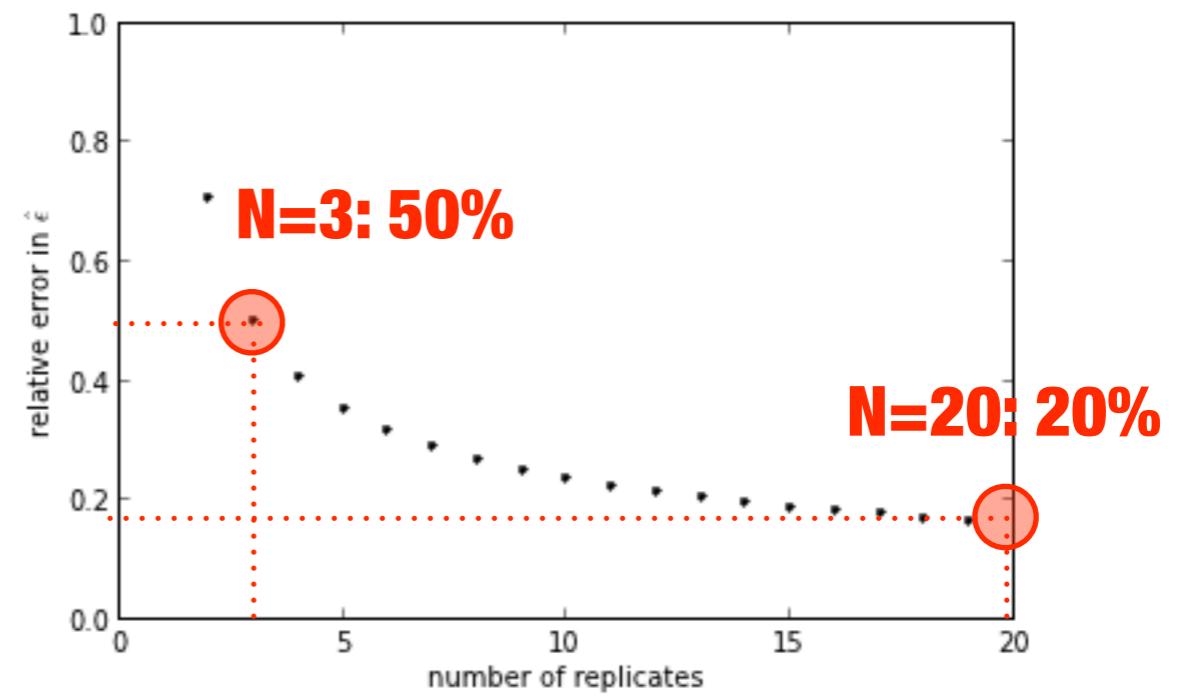
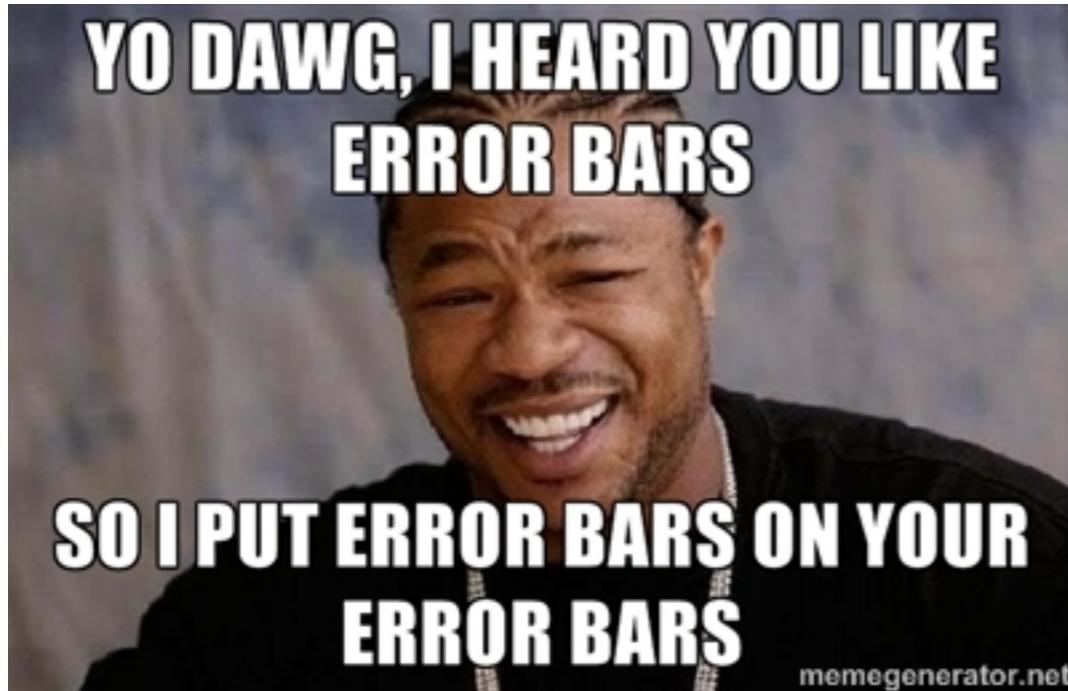
Confidence intervals [low,high] or box-and-whisker plots are also acceptable, and preferred when error is highly asymmetric (e.g. probabilities).

YES: 95% confidence interval is [4.08, 4.15]

YES: box-and-whisker plot showing quartiles



What is the **error** in the **error**?



For large N , relative error in standard error estimate is

$$\frac{\delta \hat{\epsilon}}{\hat{\epsilon}} = \frac{1}{\sqrt{2(N - 1)}}$$

For 3 replicates, the uncertainty in the error is over 50%!

Only the order of magnitude of error is known; not even a single digit is certain!

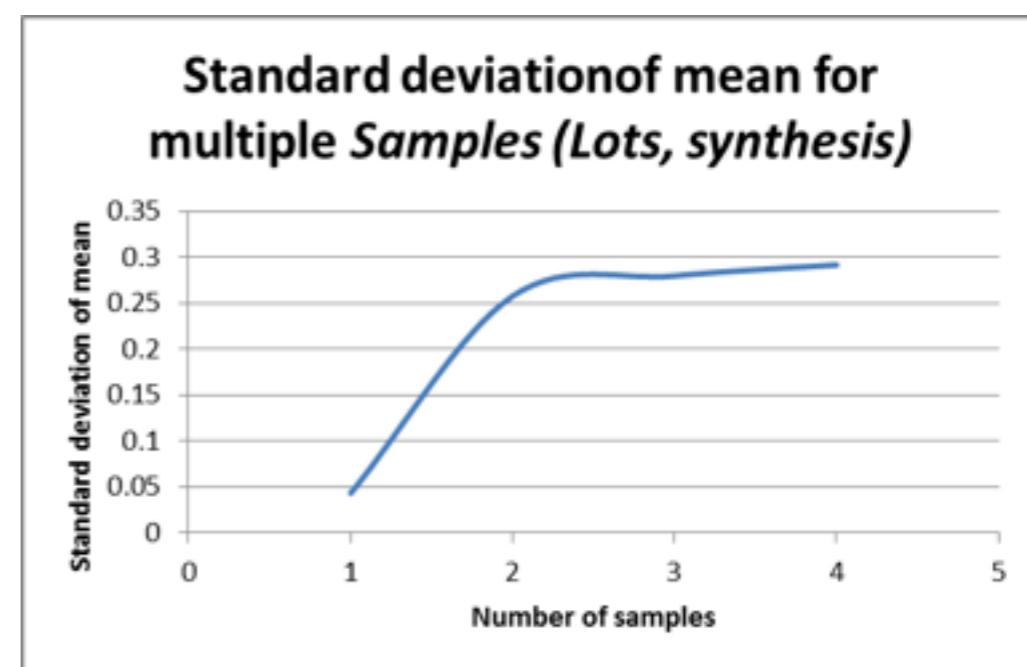
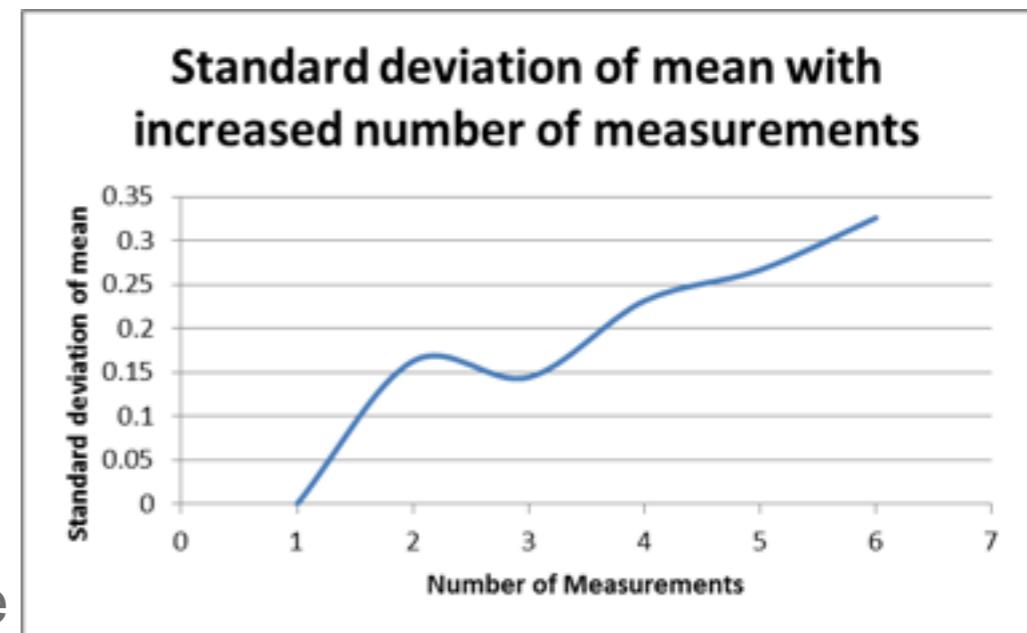
Even for 20 replicates, uncertainty is still close to within 20%!

What accuracy do we need to be useful?

- Depends on the need and use
 - Efficacy optimization
 - Factor of 2
 - HTS hit finding campaign
 - Factor of 10
 - Alerts
 - Factor of 100?
- Proper presentation of results is vital
 - “Cutoffs” and classifications can make interpretation tricky for quantitative data
 - E.g.: Issuing a hERG alert at $< 1 \mu\text{M}$
 - Propagate error into rational decision making (*a la* Optibrium)
 - E.g.” A qualitative solubility assay (i.e. soluble/insoluble DMSO/water precipitate) reported to 3 significant figures

Estimating error when none is reported

- Previous experience
- Variation in
 - Standards
 - Samples of the same compound
 - Multiple measurements of the same compound
 - Between days / equipment
- If in doubt suspect a value of 5 – 10 fold
 - Experimentalists all note for same SOP in same company
 - Different labs, instruments, technicians
 - 2X difference expected
 - 5X easy to find
 - 10X not unexpected



What is the expected error for different procedures

- From a screening veteran:
 - We are routinely using LC-MS measurement for identity check and purity of library compounds.
 - Of course there are major limitations: which method do you use, which gradient, which wavelength....
 - And even when you have settled on one method you can get two different results if the compound degrades in the solvent while it is sitting in the plate, the compound is not picked in a precise volume....
 - But it is good enough to check if there is *major* degradation going on.
 - To be more precise you would need to use NMR but that's not feasible for a *classical library compound* – only after re-synthesis.

What is the expected error

- Purity: Corporate archive: 10% of the compounds were found to have degraded to < 75% purity (Kramer and Lewis, 2012)
 - IMHO this might be an optimistic number for many companies
 - At one time, purity was defined at 70% in some 80% in other companies
- Pipetting of serial dilutions: The *mixing* of the compound may be done entirely different using different methods resulting in 10 fold shifts of IC50 when using one method over the other. (Well informed screening veteran)
 - Angle of the pipette affects mixing



Stouch

CADD GRC July 2013

Modeling experimental error can be a valuable exercise: Biomek dilution series vs Echo acoustic dispensing

In the Pipeline

[« Aveo Gets Bad News on Tivozanib | Main | The Medical Periodic Table »](#)

May 3, 2013

Drug Assay Numbers, All Over the Place

Posted by Derek

There's a [truly disturbing paper](#) out in PLoS ONE with potential implications for a lot of assay data out there in the literature. The authors are looking at the results of biochemical assays as a function of how the compounds are dispensed in them, pipet tip versus **acoustic**, which is the sort of idea that some people might roll their eyes at. But people who've actually done a lot of biological assays may well feel a chill at the thought, because this is just the sort of you're-kidding variable that can make a big difference.

http://pipeline.corante.com/archives/2013/05/03/drug_assay_numbers_all_over_the_place.php

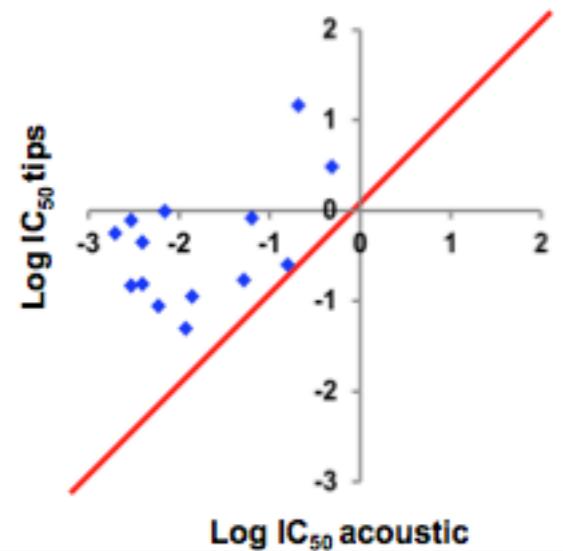
OPEN  ACCESS Freely available online



Dispensing Processes Impact Apparent Biological Activity as Determined by Computational and Statistical Analyses

Sean Ekins^{1*}, Joe Olechno², Antony J. Williams³

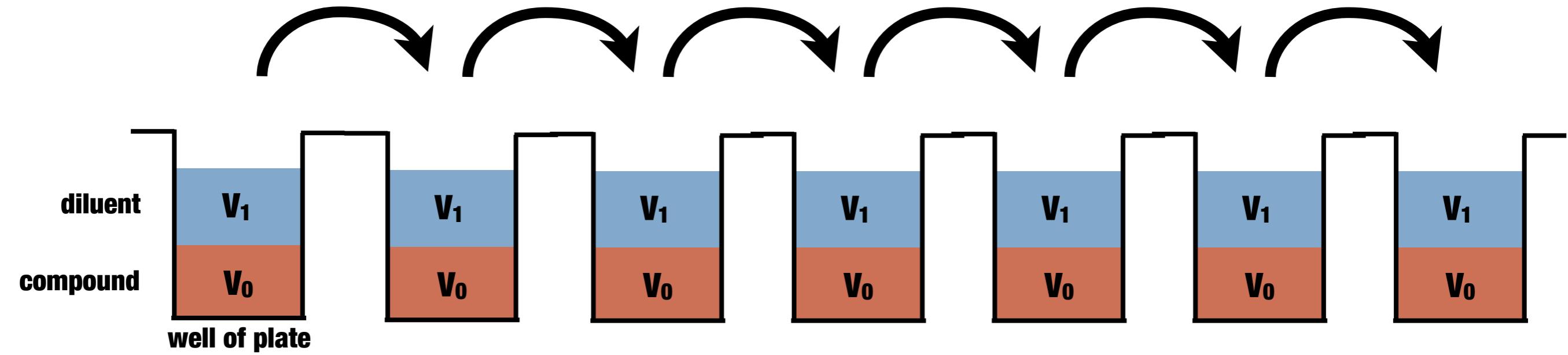
¹ Collaborations in Chemistry, Fuquay-Varina, North Carolina, United States of America, ² Labcyte Inc., Sunnyvale, California, United States of America, ³ Royal Society of Chemistry, Wake Forest, North Carolina, United States of America



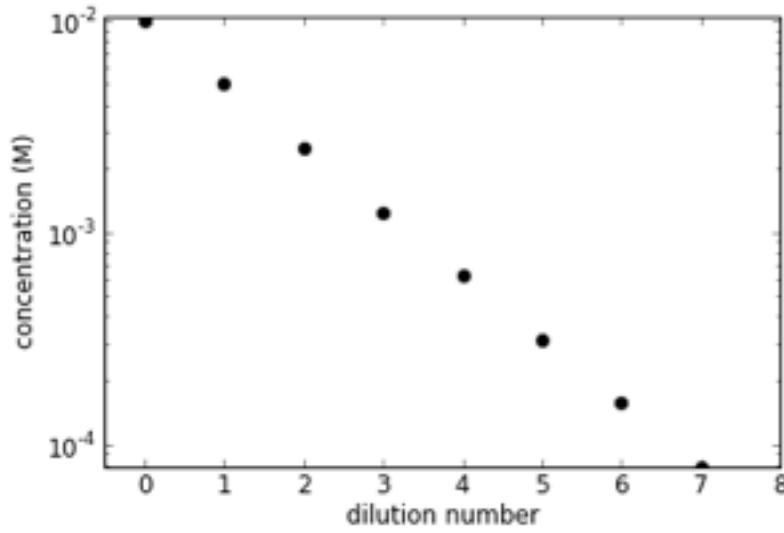
Elkins et al. PLoS One 8:e62325, 2013.

Is this just a result of pipetting accuracy differences?
Let's model the experiment to find out.

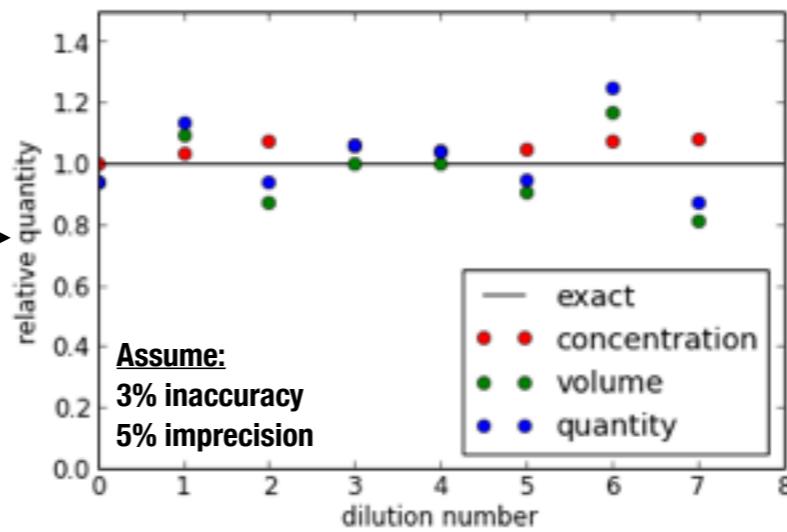
Modeling experimental error can be a valuable exercise: Preparing a dilution series via liquid-handling robot



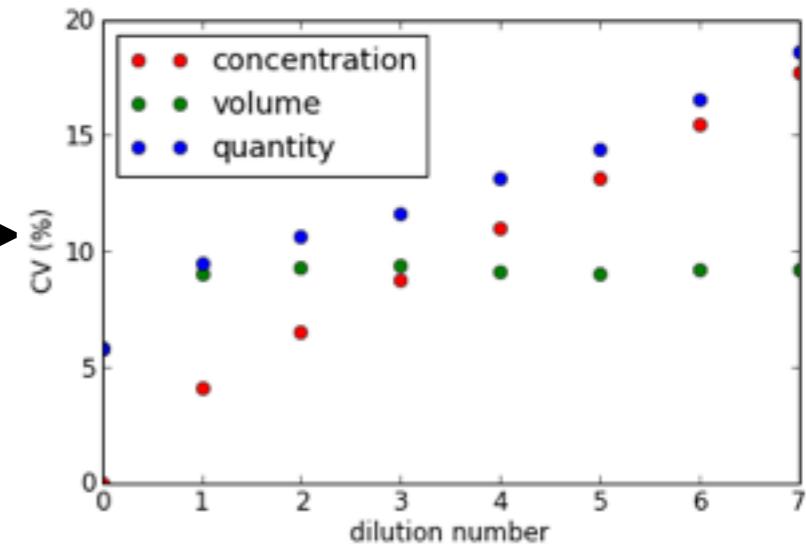
ideal concentrations



error for a single realization



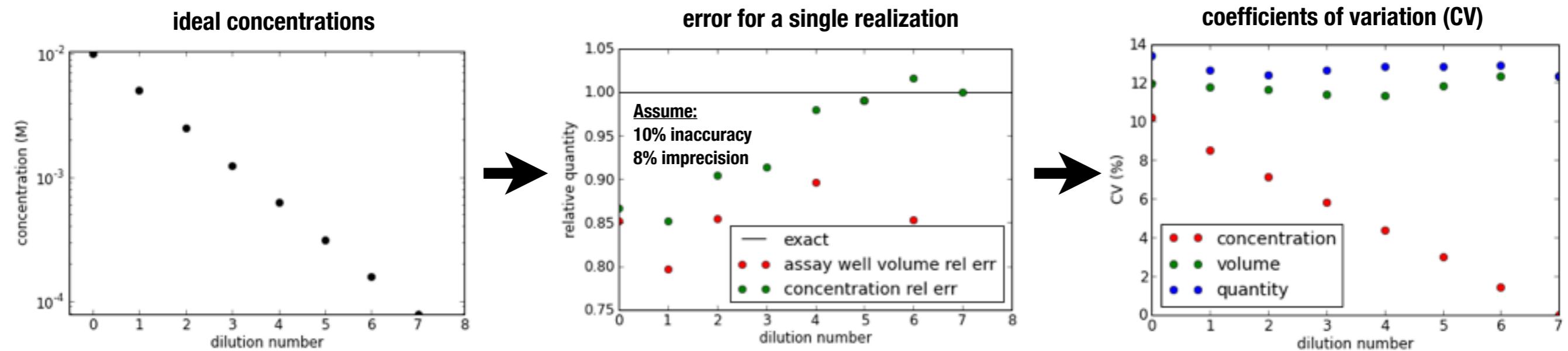
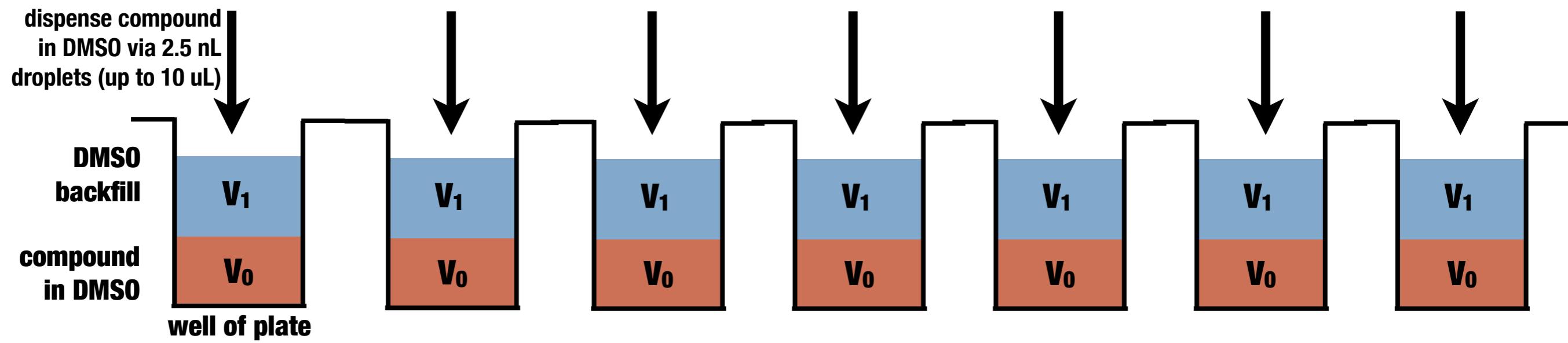
coefficients of variation (CV)





**An Inconvenient Truth:
It's hard to make a good dilution series**

Modeling experimental error can be a valuable exercise: Preparing a dilution series via acoustic dispensing

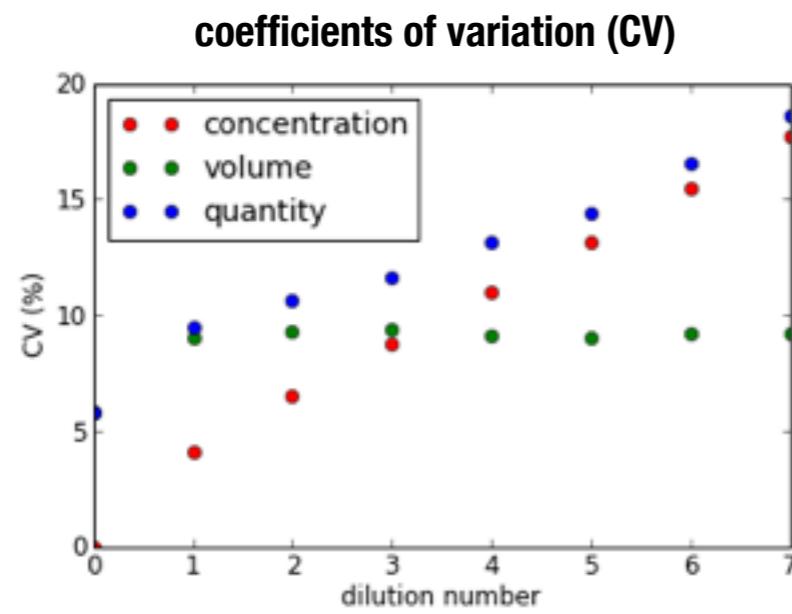


http://www.labcyte.com/sites/default/files/support_docs/Echo%205XX%20Specifications.pdf

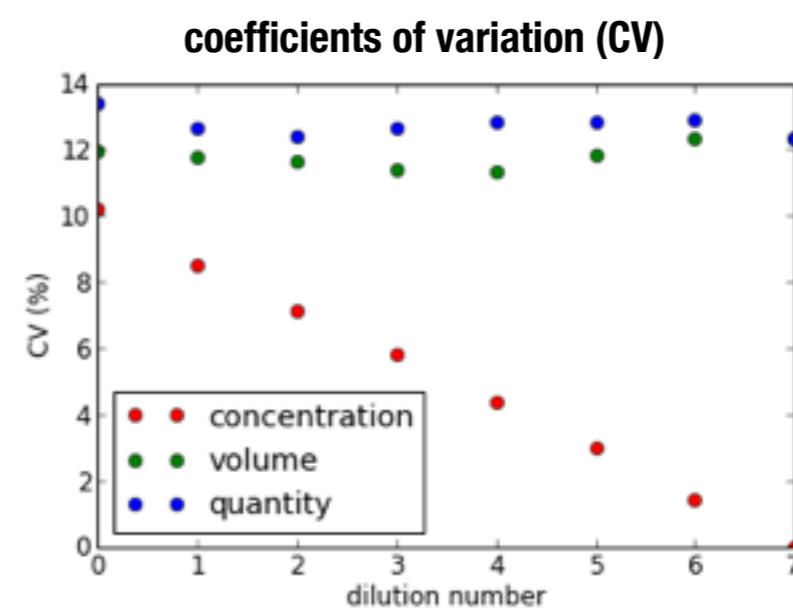


LabCyte Echo

Modeling experimental error can be a valuable exercise: Comparison of tip-based and acoustic dispensing



tip-based



acoustic

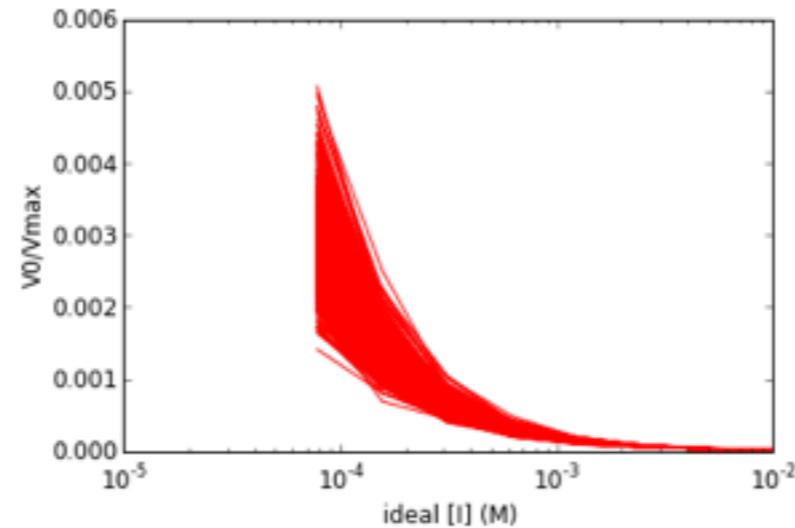
Modeling experimental error can be a valuable exercise: Assay model

↓
Dispense into assay plate:
2 uL compound dilution
10 uL enzyme assay mix



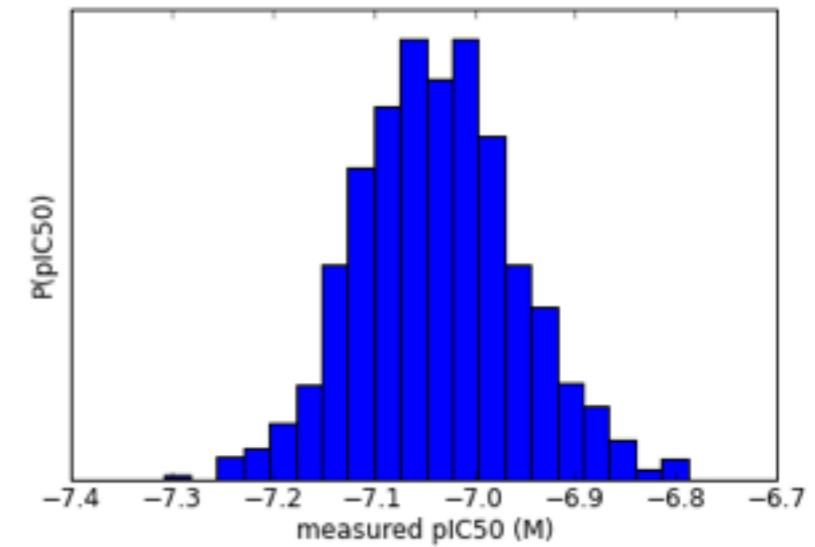
Assume for EphB4 assay:
Michaelis-Menten kinetics
competitive inhibition
[ATP] = 4 uM
[EphB4] = 6 uM
Km(ATP) = 1.71 uM

measure initial reaction velocity



fit it to get IC50

distribution of fit pIC50s



Modeling experimental error can be a valuable exercise: Biomek dilution series vs Echo acoustic dispensing

In the Pipeline

[« Aveo Gets Bad News on Tivozanib | Main | The Medical Periodic Table »](#)

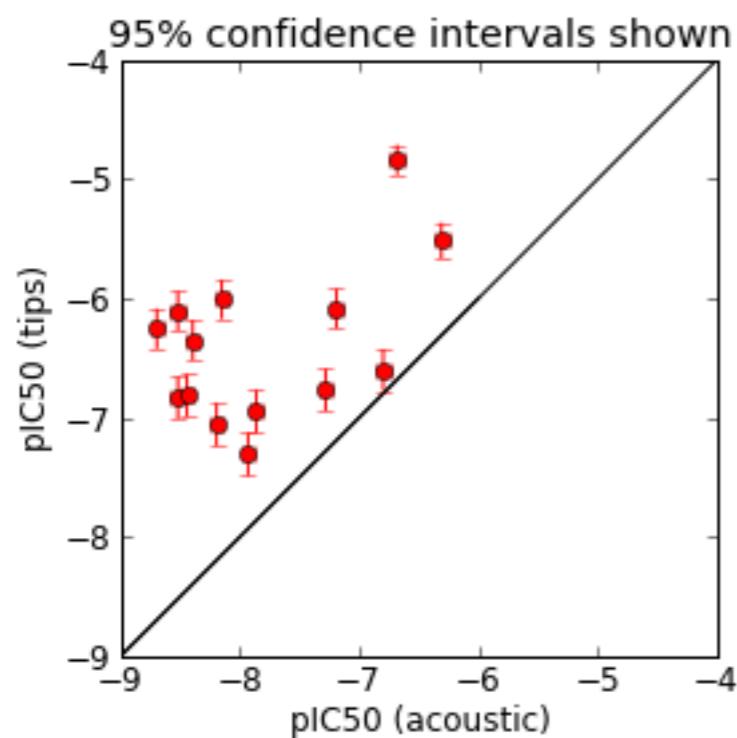
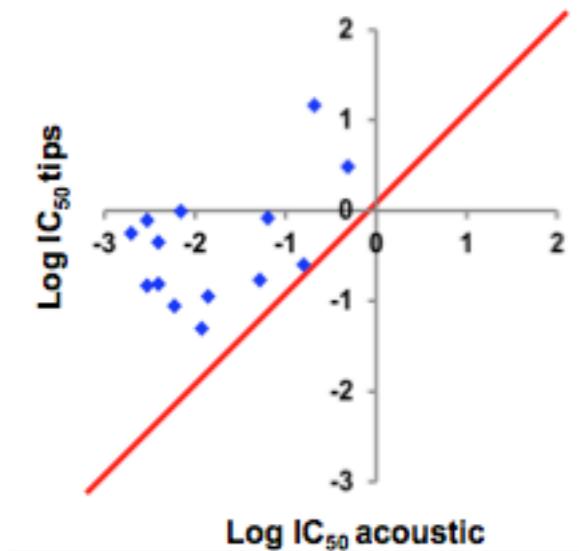
May 3, 2013

Drug Assay Numbers, All Over the Place

Posted by Derek

There's a [truly disturbing paper](#) out in PLoS ONE with potential implications for a lot of assay data out there in the literature. The authors are looking at the results of biochemical assays as a function of how the compounds are dispensed in them, pipet tip versus **acoustic**, which is the sort of idea that some people might roll their eyes at. But people who've actually done a lot of biological assays may well feel a chill at the thought, because this is just the sort of you're-kidding variable that can make a big difference.

http://pipeline.corante.com/archives/2013/05/03/drug_assay_numbers_all_over_the_place.php



Elkins et al. PLoS One 8:e62325, 2013.

OK, that's not sufficient to explain the discrepancy...
What are we missing?

Modeling experimental error can be a valuable exercise: Dilution is a big problem with liquid-based pipetting

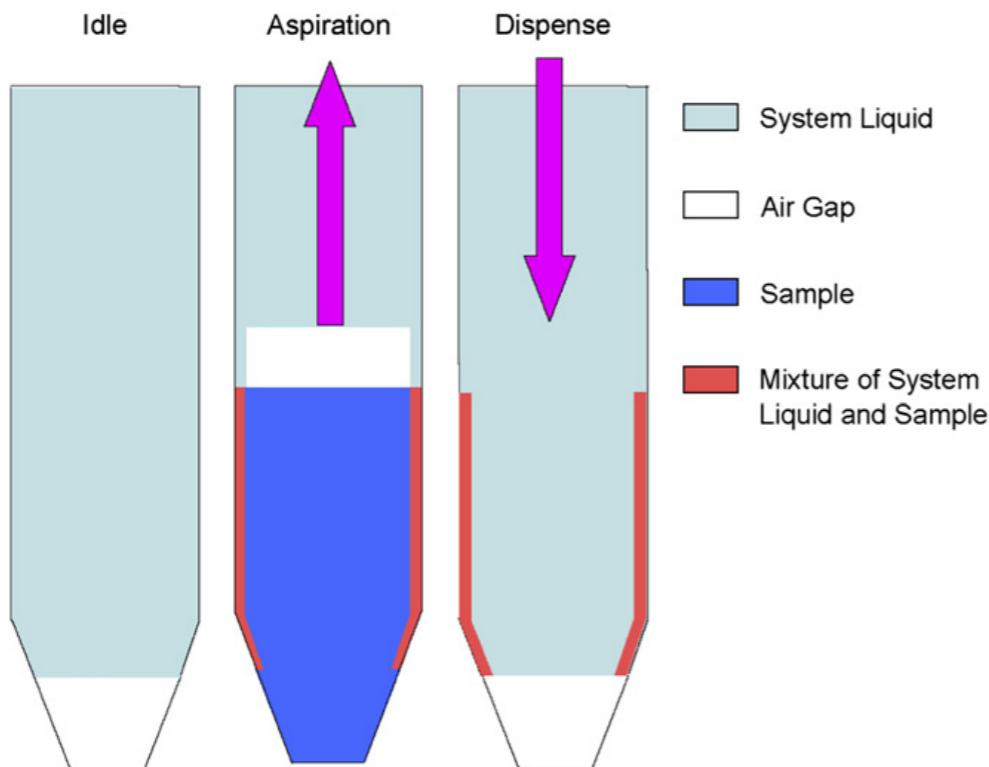


Figure 1. Schematic diagram of the dilution effect.

Table I. Comparison between the MVS method and the gravimetry method—Tecan ALH using water liquid class

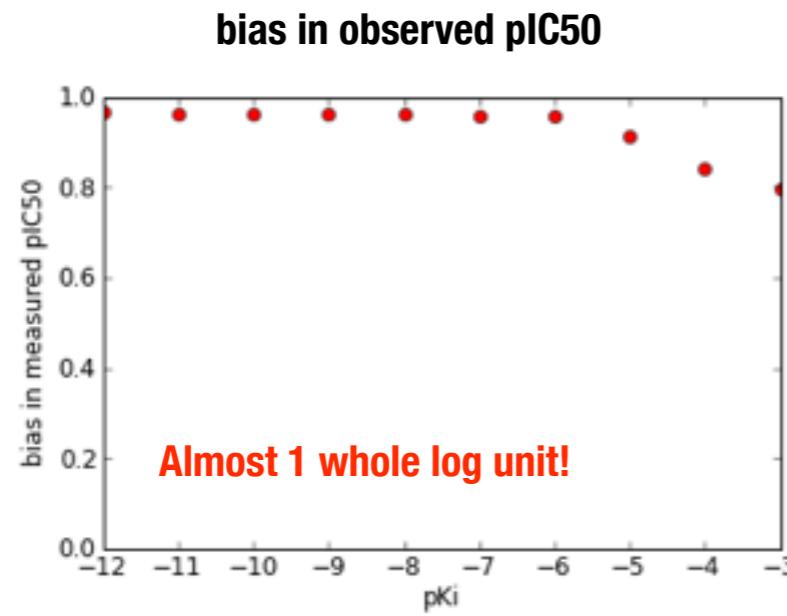
| | MVS (μL) | Gravimetry (μL) | |
|-------------------------------|-----------------------|------------------------------|-------|
| Target volume | 20 | 200 | 20 |
| Mean volume (μL) | 18.74 | 190.08 | 20.15 |
| Inaccuracy (%) | -6.30 | -4.96 | 0.75 |
| StDev | 0.22 | 1.74 | 0.6 |
| CV ($n = 96$) (%) | 1.17 | 0.92 | 1.94 |
| | | | 0.59 |

These solutions were delivered by the Tecan ALH using an aspirate/dispense protocol based on a water liquid class. Inaccuracy corresponds to $100 \times (\text{Observed volume} - \text{Target volume})/\text{Target volume}$. Coefficient of variation (CV) corresponds to $100 \times \text{StDev}/\text{mean}$, where StDev is the standard deviation.

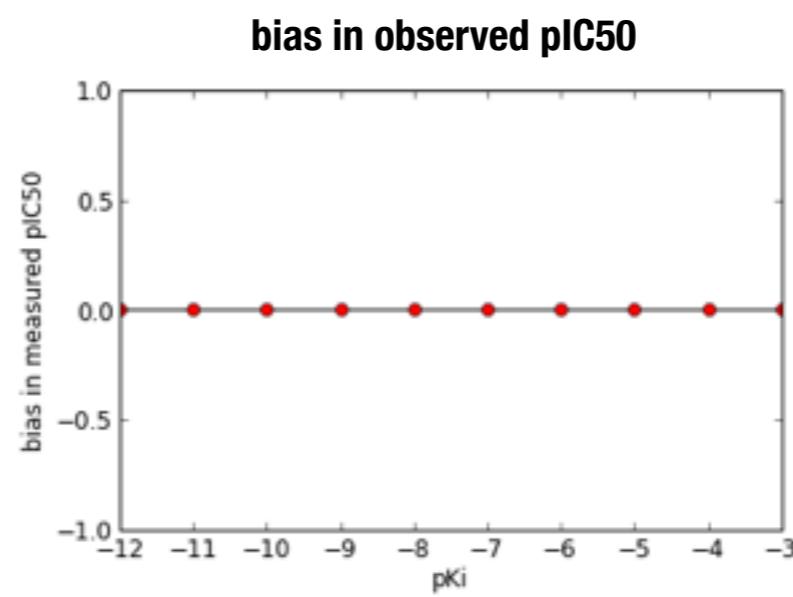


Gu and Deng. JALA 12:355, 2007 (BMS)
Dong, Ouyang, Liu, and Jemal. JALA 11:60, 2006. (BMS)

Modeling experimental error can be a valuable exercise: Comparison of tip-based and acoustic dispensing



tip-based



acoustic

Modeling experimental error can be a valuable exercise: Biomek dilution series vs Echo acoustic dispensing

In the Pipeline

[« Aveo Gets Bad News on Tivozanib | Main | The Medical Periodic Table »](#)

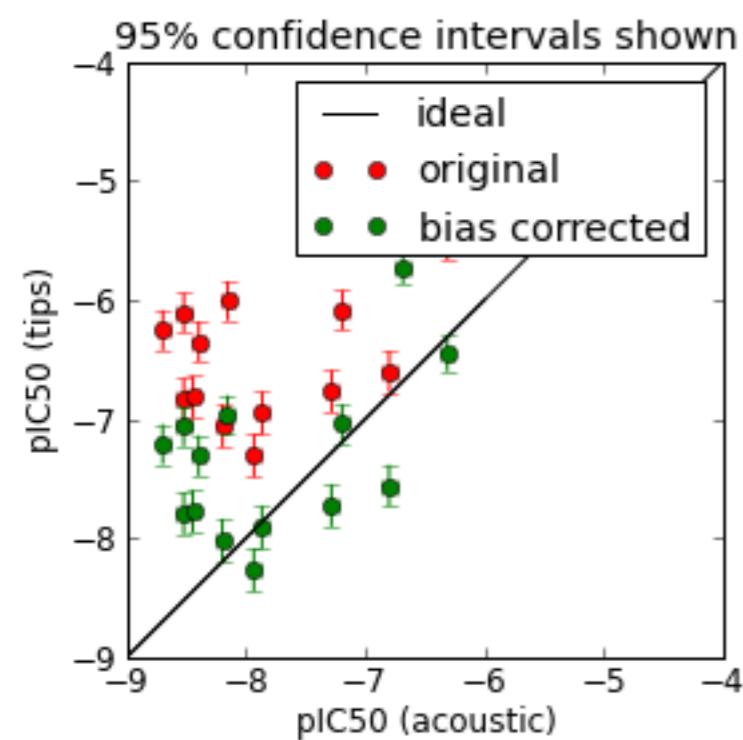
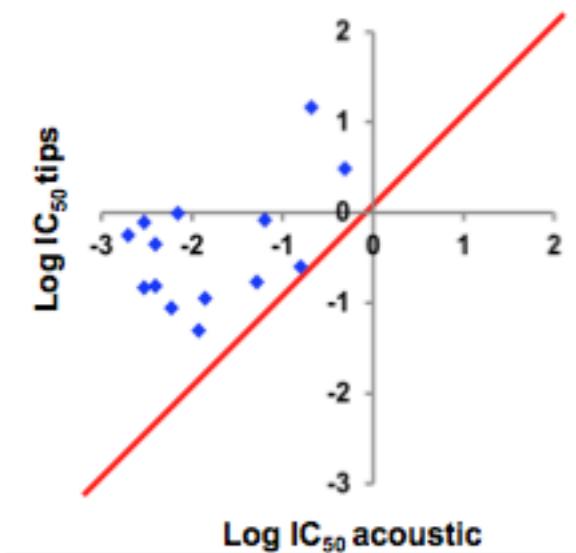
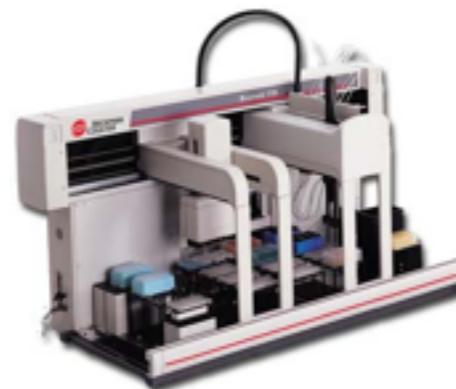
May 3, 2013

Drug Assay Numbers, All Over the Place

Posted by Derek

There's a [truly disturbing paper](#) out in PLoS ONE with potential implications for a lot of assay data out there in the literature. The authors are looking at the results of biochemical assays as a function of how the compounds are dispensed in them, pipet tip versus **acoustic**, which is the sort of idea that some people might roll their eyes at. But people who've actually done a lot of biological assays may well feel a chill at the thought, because this is just the sort of you're-kidding variable that can make a big difference.

http://pipeline.corante.com/archives/2013/05/03/drug_assay_numbers_all_over_the_place.php



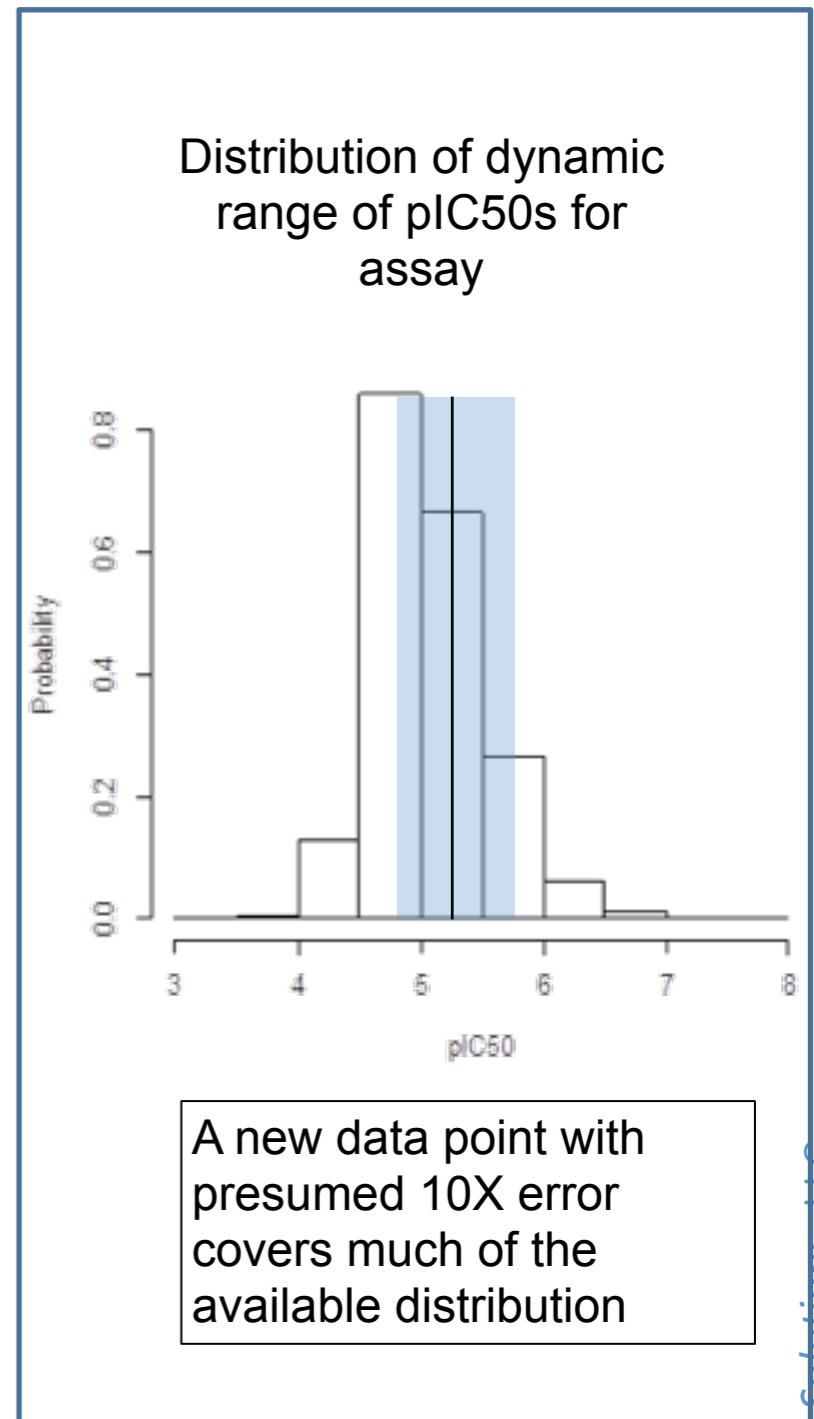
Elkins et al. PLoS One 8:e62325, 2013.

How bad can things really be? (Catastrophic failures)

- Worst case: Misled: The data is useless but does not appear to be so
 - Actual error is greater than the data range
 - Actual error is greater than the needed precision
 - Data is reported with the wrong precision
 - Data is (allowed to be) used incorrectly
- Next worst: Bereft: The data is useless

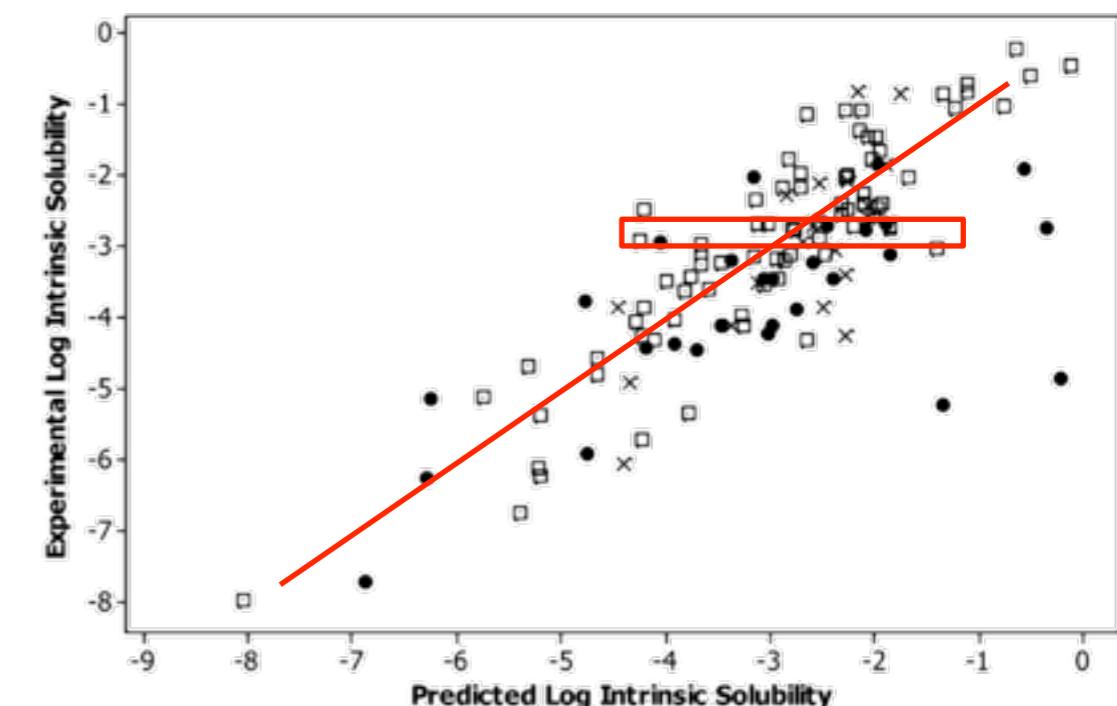
How bad can things really be? (Catastrophic failures)

- Worst case: Misled: The data is useless but does not appear to be so
 - Actual error is greater than the data range
 - Actual error is greater than the needed precision
 - Data is reported with the wrong precision
 - Data is (allowed to be) used incorrectly
- Next worst: Bereft: The data is useless
- Factor of 10 between sites is not uncommon
 - Considering the limited range of much of Pharma data a factor of 10 might swamp the results
 - i.e. The error might be a substantial portion of the dynamic range of the data
 - In particular, the data might not be useful for a desired use
 - E.g. Optimization vs. alert



How safe is it to use collated literature data?

- In 2 words, not very
- Depends on the property as well as how the data and model are used
- Assay conditions can vary widely
- It rarely happens that details of the experiments and standards or common molecules are available
 - Are there common molecules that can be used for calibration?
- Possibly use for approximate, general models of low expected precision
 - Some modest successes can be had
- Expect a factor of 10 error



M. Hewitt, et al 2009, J. Chem. Inf. Model, 49, 2572–2587

What if error was not reported or measured? How reliable is public pKi data?

Journal of
**Medicinal
Chemistry**

Article

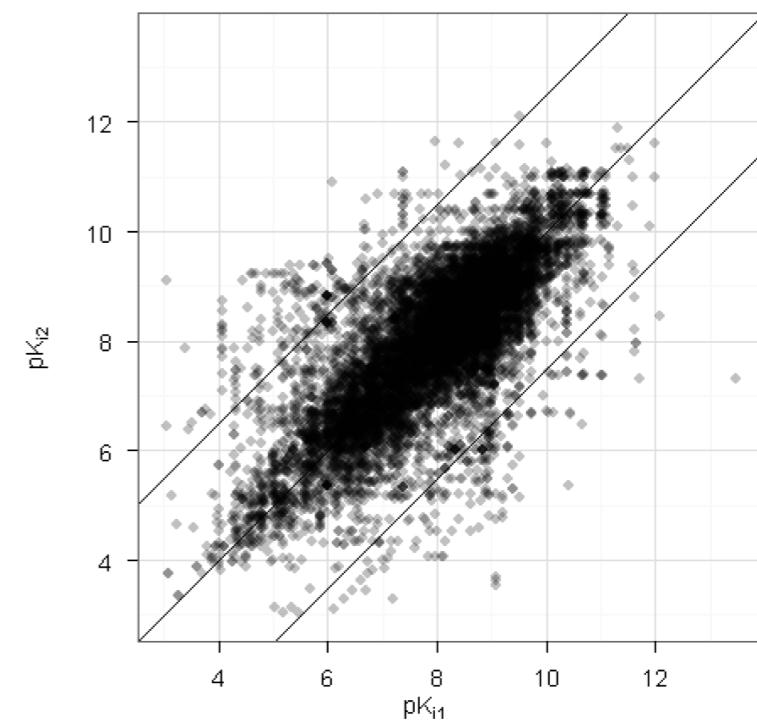
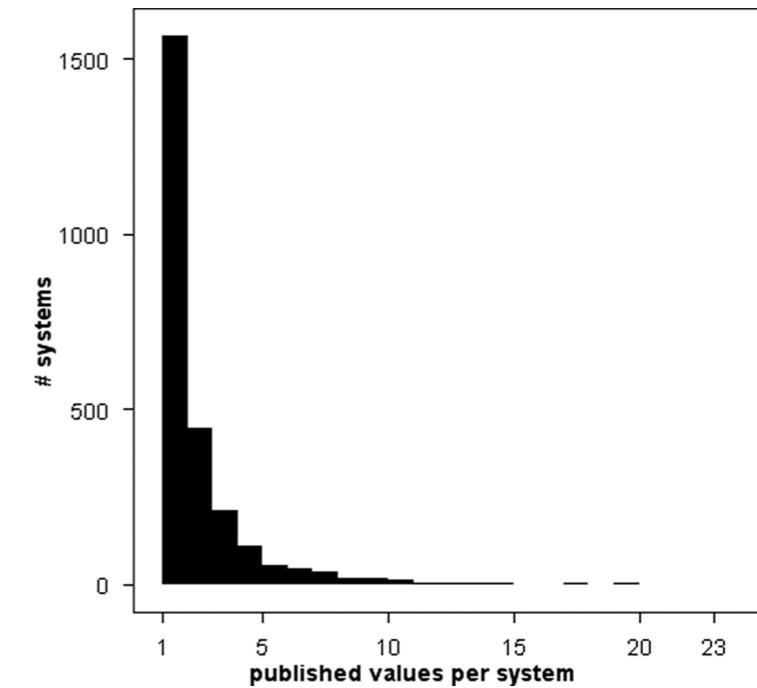
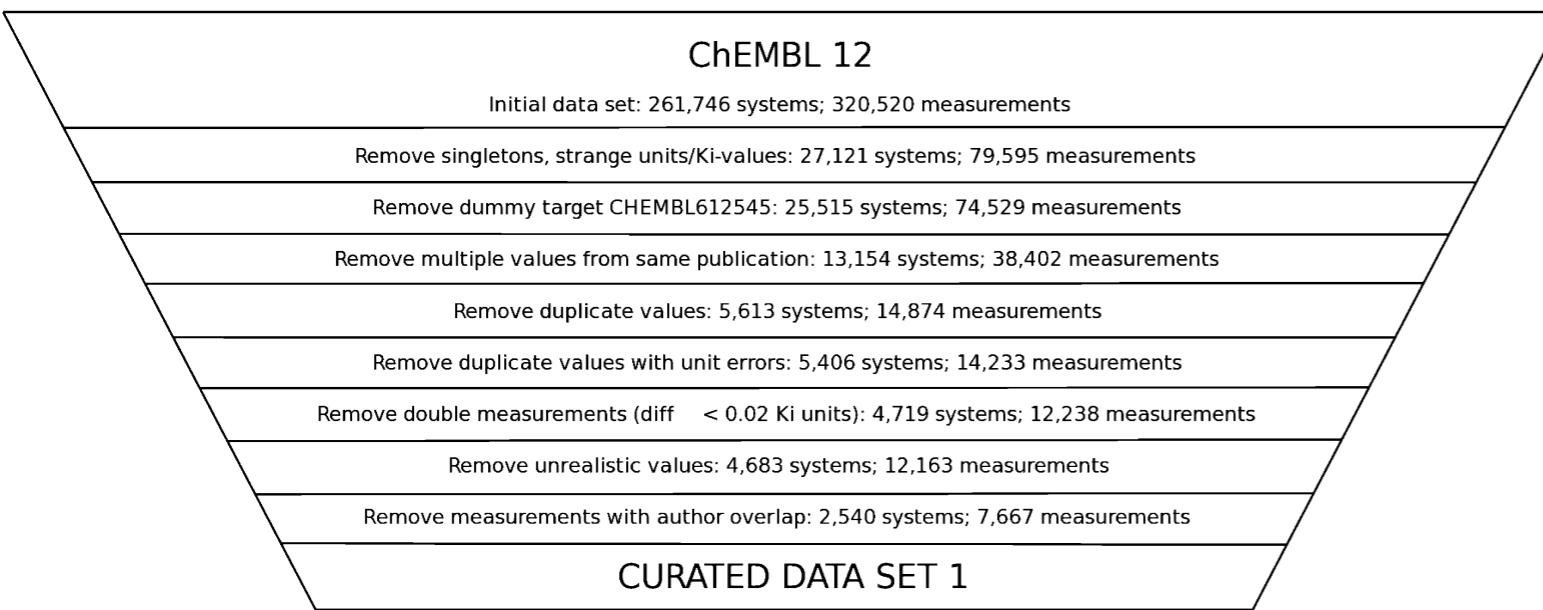
pubs.acs.org/jmc

The Experimental Uncertainty of Heterogeneous Public K_i Data

Christian Kramer,^{*,†} Tuomo Kalliokoski,^{*,†} Peter Gedeck, and Anna Vulpetti

Novartis Institutes for BioMedical Research, Novartis Pharma AG, Forum 1, Novartis Campus, CH-4056 Basel, Switzerland

Data curation

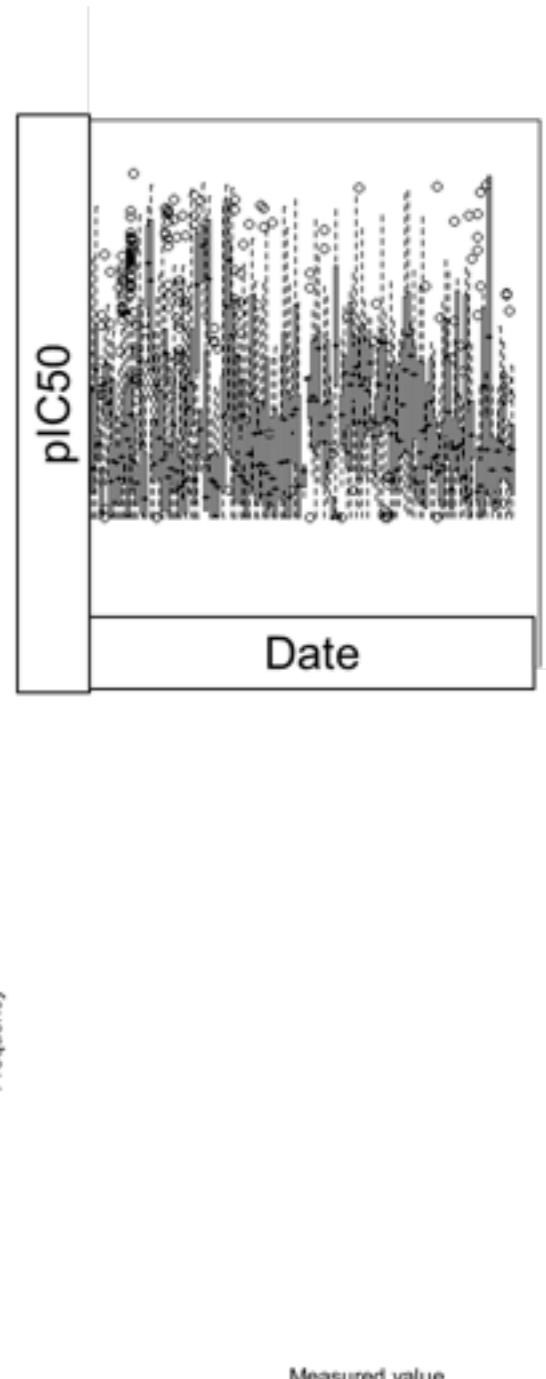


root mean square error (RMSE) ~ 0.54 pKi units ~ 0.75 kcal/mol @ 25 C

Christian Kramer et al. J Med Chem 2012.

Spotting "place holder" values

- (0.0 is really a bad default for most Pharma data)
- Depends on the data form and range
- Suspicious, recurring value – particularly for continuous data
- Values way outside the norm (9999.9999)
- Blank values (depending on how they will be interpreted – set to 0.0?)
- Numbers that occur too often and/or are too ideal
 - Ex: 4.000 as the most frequent response whereas there were few 4 sig fig duplicates otherwise



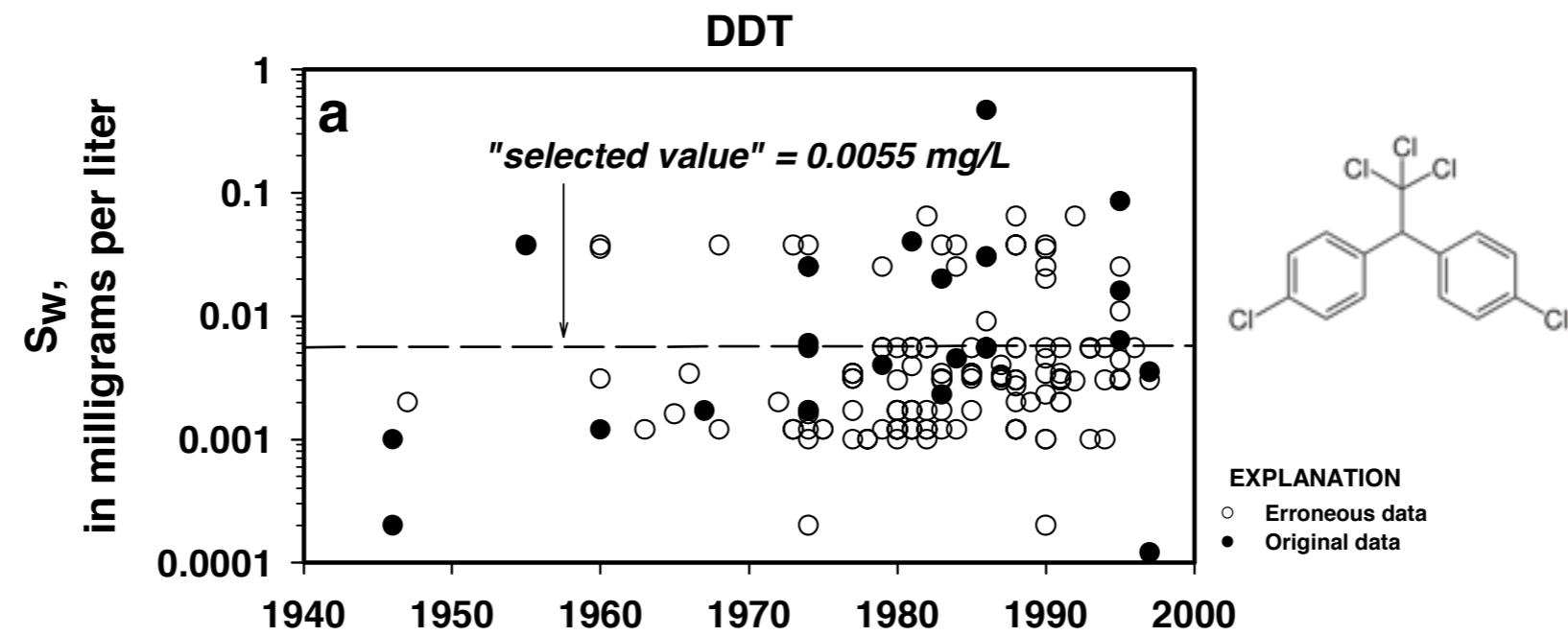
Spotting "place holder" values

- The danger of “0”
- Is there metadata suggesting that an experiment was not done or interrupted or somehow failed?
 - Example:
 - Mean value for logD of 1.85 (the automatically provided value)
 - 4 measurements
 - 3.7, 0.0, 4.0, 0.0
 - Off in the corner was a free form note that noted for the 0.0 values
 - ~"Insufficient material to complete experiment"
- If one relied on the default average value, the value would be 2 log units in error

Pitfalls with literature data and public databases

The accurate determination of an organic contaminant's physico-chemical properties is essential for predicting its environmental impact and fate. Approximately 700 publications (1944–2001) were reviewed and all known aqueous solubilities (S_w) and octanol-water partition coefficients (K_{ow}) for the organochlorine pesticide, DDT, and its persistent metabolite, DDE were compiled and examined. Two problems are evident with the available database: 1) egregious errors in reporting data and references, and 2) poor data quality and/or inadequate documentation of procedures. The published literature (particularly the collative literature such as compilation articles and handbooks) is characterized by a preponderance of unnecessary data duplication. Numerous data and citation errors are also present in the literature. The percentage of original S_w and K_{ow} data in compilations has decreased with time, and in the most recent publications (1994–97) it composes only 6–26 percent of the reported data. The variability of original DDT/DDE S_w and K_{ow} data spans 2–4 orders of magnitude, and there is little indication that the uncertainty in these properties has declined over the last 5 decades. A criteria-based evaluation of DDT/DDE S_w and K_{ow} data sources shows that 95–100 percent of the database literature is of poor or unevaluatable quality. The accuracy and reliability of the vast majority of the data are unknown due to inadequate documentation of the methods of determination used by the authors. [For example, estimates of precision have been reported for only 20 percent of experimental S_w data and 10 percent of experimental K_{ow} data.] Computational methods for estimating these parameters have been increasingly substituted for direct or indirect experimental determination despite the fact that the data used for model development and validation may be of unknown reliability. Because of the prevalence of errors, the lack of methodological documentation, and unsatisfactory data quality, the reliability of the DDT/ DDE S_w and K_{ow} database is questionable. The nature and extent of the errors documented in this study are probably indicative of a more general problem in the literature of hydrophobic organic compounds. Under these circumstances, estimation of critical environmental parameters on the basis of S_w and K_{ow} (for example, bioconcentration factors, equilibrium partition coefficients) is inadvisable because it will likely lead to incorrect environmental risk assessments. The current state of the database indicates that much greater efforts are needed to: 1) halt the proliferation of erroneous data and references, 2) initiate a coordinated program to develop improved methods of property determination, 3) establish and maintain consistent reporting requirements for physico-chemical property data, and 4) create a mechanism for archiving reliable data for widespread use in the scientific/regulatory community.

The literature is filled with erroneous data



Data has a habit of being re/misreported

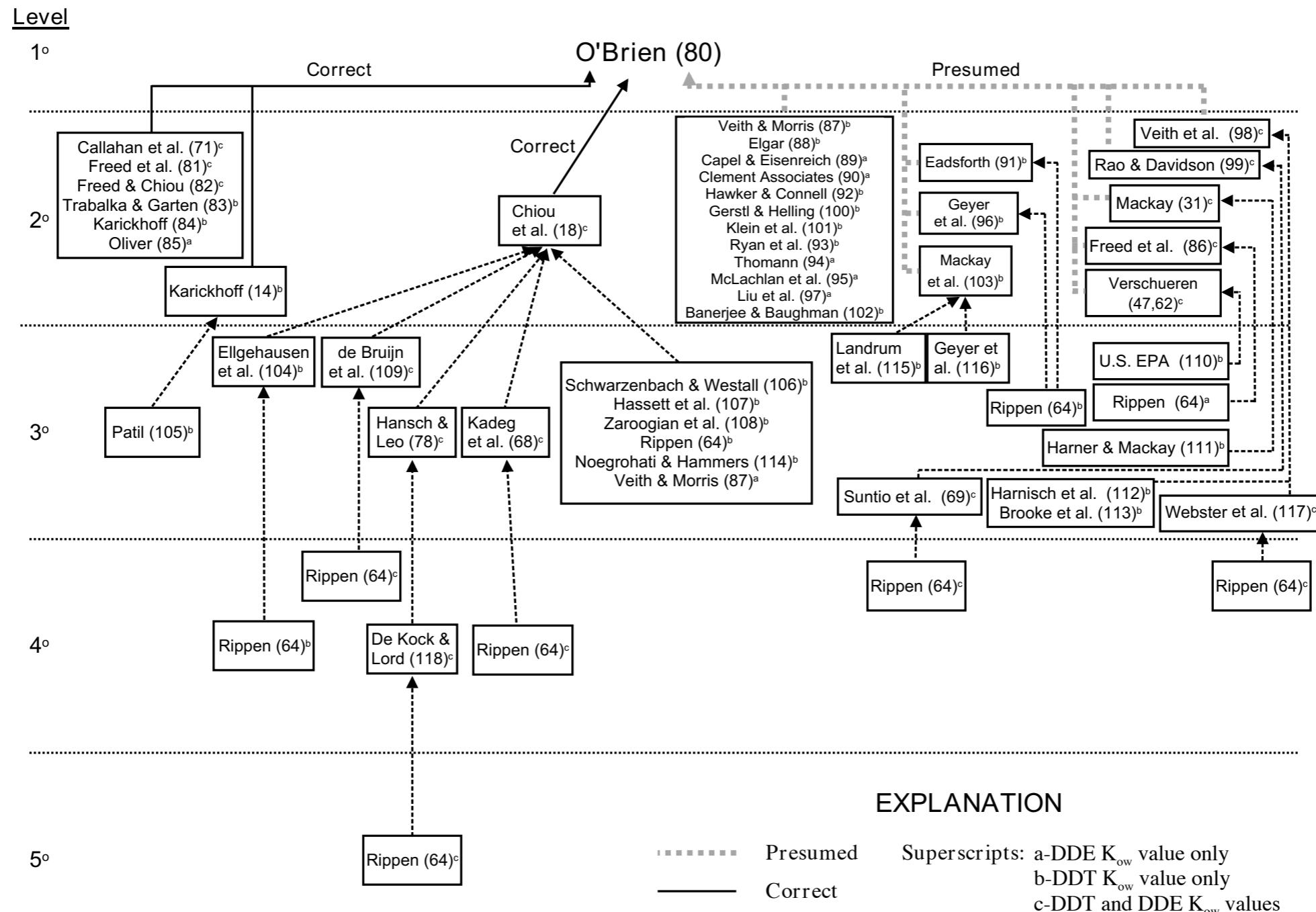


Figure 1. Reference tree showing original O'Brien (80) publication and secondary, tertiary and higher references as an example of multi-level referencing.

USGS Water-Resources Investigations Report 01-4201, 2001.

What if error was not reported or measured? Reported measurements can vary over huge ranges

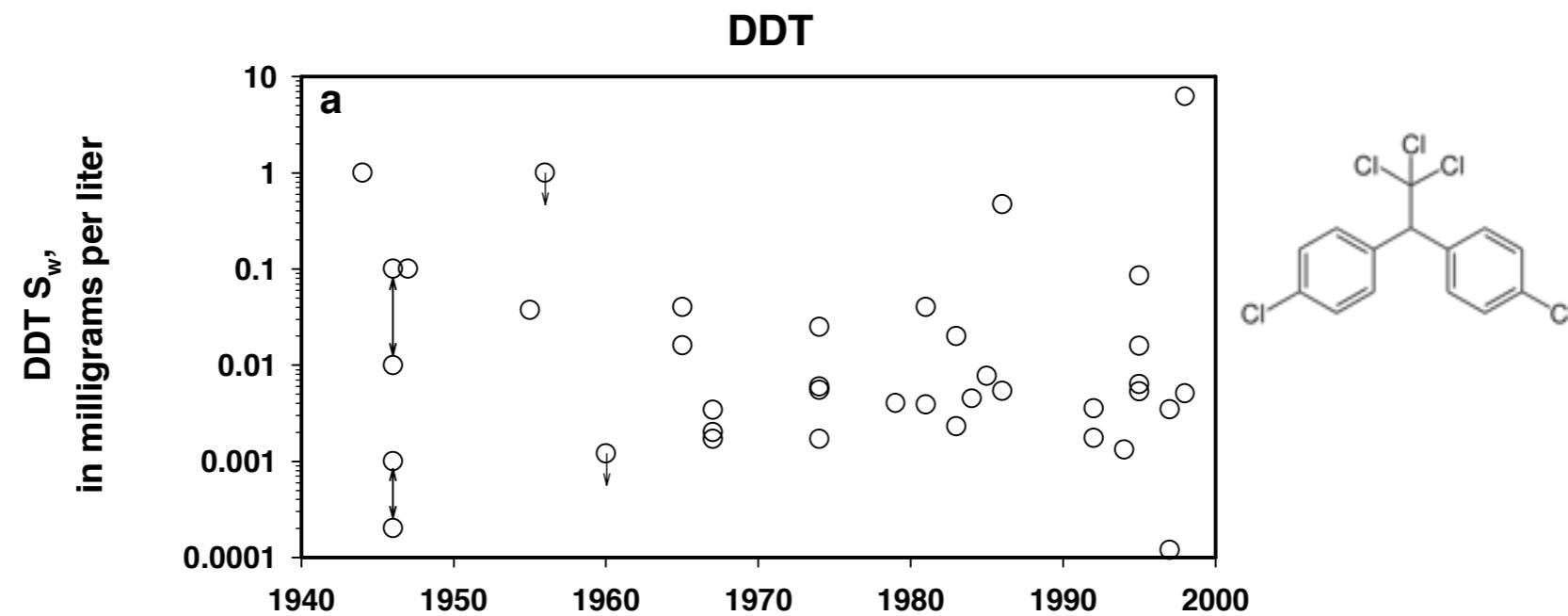


Table 6. Basic Statistics for Uncensored Original S_w and K_{ow} DDT/DDE Data^a

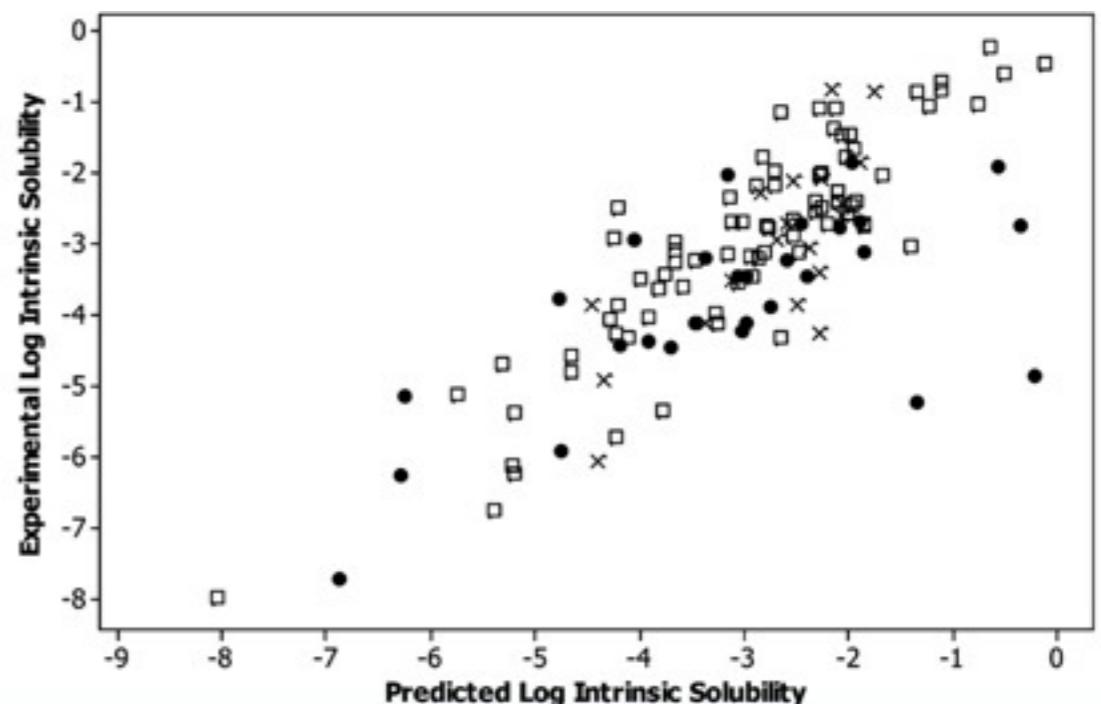
| <u>Statistic</u> | DDT | DDE |
|--------------------|---|------------|
| | <u>Aqueous solubility, S_w (milligrams per liter)</u> | |
| Mean | 0.23 | 0.12 |
| Standard Deviation | 1.00 | 0.32 |

Noisy bad data – better than none?

- When is a lot of noisy data worse than much-less but high quality data?
 - Always
 - Always, *unless*, there is a reliable core of high quality data to which the lesser quality data can be compared
- When is NO measurement actually a better option?
 - Often
 - Example: Much high and mid throughput ADME data is not of a precision to help in optimization efforts

When is lots of noisy data worse than less higher quality data?

- Depends on the need. Will the data be used for:
 - Alerts
 - Crude general indicator model
 - Optimization
- Perhaps useful when large dynamic range is required and is more important than precision
 - E.g. for crude solubility models



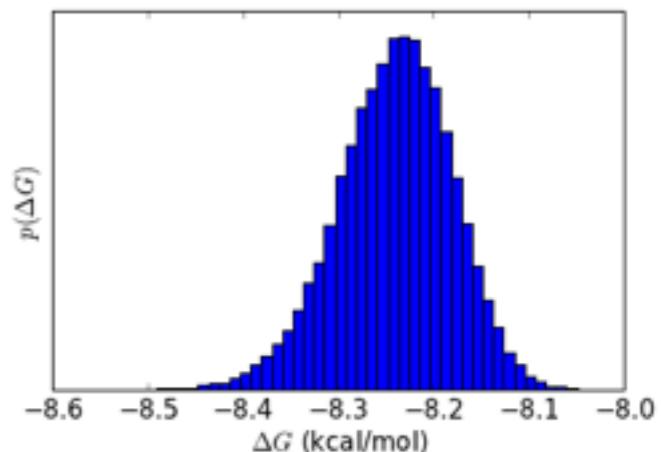
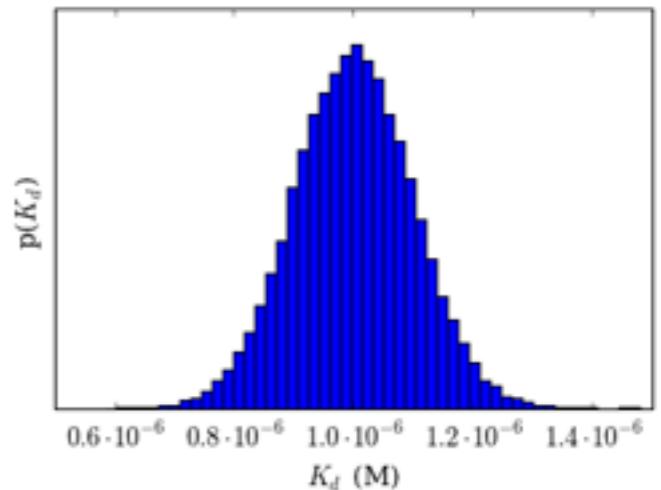
Propagation of (t)error

**Suppose we make a measurement of the K_d and find 1.0 ± 0.1 uM.
(Assume the error in K_d is Gaussian)**

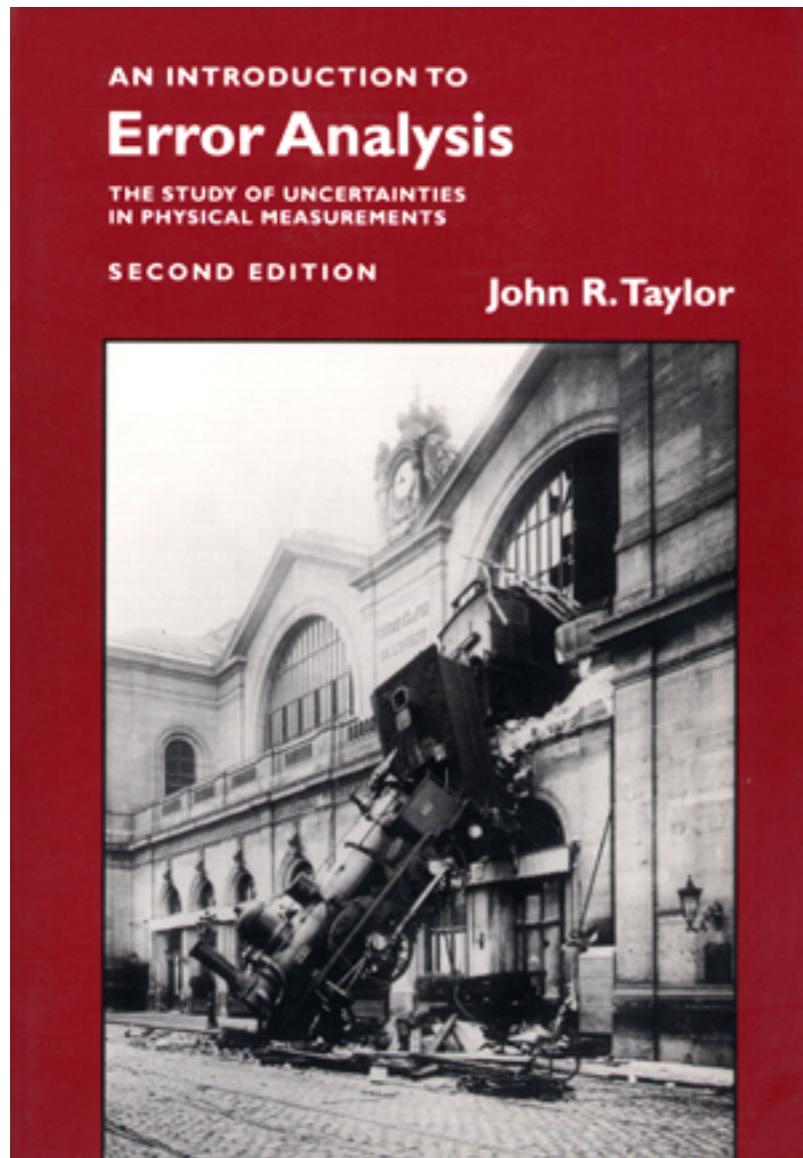
Suppose we want to compute the free energy of binding ΔG from the K_d

$$\Delta G_{\text{bind}} = k_B T \ln(K_d C_0)$$

How can we propagate the error δK_d into ΔG ?



Propagation of (t)error: What the book says



Adding two measurements with error:

$$z = x + y$$

$$\delta^2 z = \delta^2 x + \delta^2 y + \delta x \delta y$$

Multiplying two measurements with error:

$$z = x \cdot y$$

$$\delta^2 z = x^2 \delta^2 y + y^2 \delta^2 x + xy \delta x \delta y$$

Dividing two measurements with error:

$$z = x/y$$

$$\delta^2 z = z^2 \left[\left(\frac{\delta x}{x} \right)^2 + \left(\frac{\delta y}{y} \right)^2 + \left(\frac{\delta x \delta y}{xy} \right) \right]$$

**Where does this stuff come from?
Why is it so complicated?
Is there something easier?**

Propagation of (t)error: Where does it come from?

Suppose we have a function of a measurement or random variable $y = f(x)$

We make a measurement with some uncertainty $x_0 \pm \delta x$

We can describe the behavior in the vicinity of the measurement using a first-order Taylor series expansion:

$$\Delta f \equiv f(x_0 + \Delta x) - f(x_0) = \left. \frac{\partial f}{\partial x} \right|_{x_0} \Delta x + \mathcal{O}(\Delta x^2)$$

The squared uncertainty is then given by its variance, omitting all but the first-order term:

$$\delta^2 f = E[\Delta f^2] = E \left[\left(\left. \frac{\partial f}{\partial x} \right|_{x_0} \Delta x \right)^2 \right] = \left(\left. \frac{\partial f}{\partial x} \right|_{x_0} \right)^2 \delta^2 x$$

The uncertainty (standard error) is just the square root of the squared uncertainty:

$$\delta f = \left| \left. \frac{\partial f}{\partial x} \right|_{x_0} \right| \delta x$$

Propagation of (t)error: Converting affinity to free energy.

**Suppose we make a measurement of the K_d and find $1.0 \pm 0.1 \text{ uM}$.
(Assume the error in K_d is Gaussian)**

Suppose we want to compute the free energy of binding ΔG from the K_d

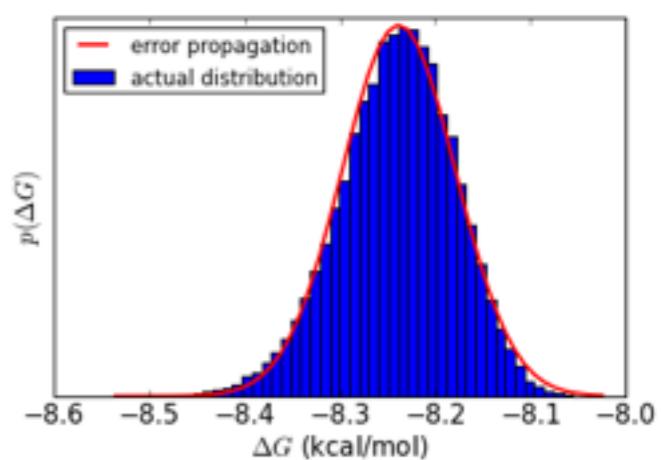
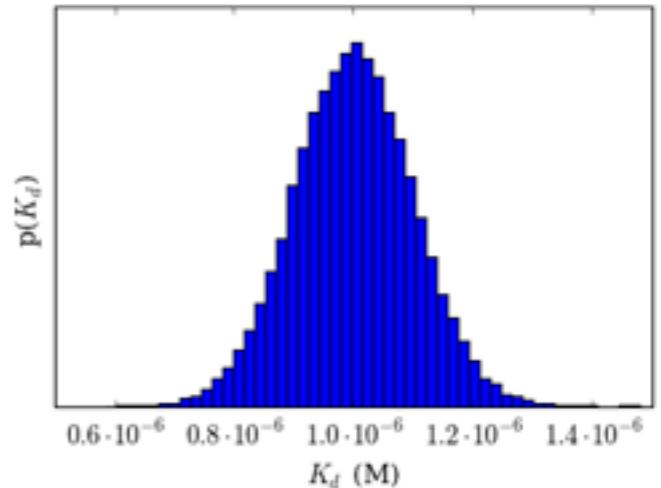
$$\Delta G_{\text{bind}} = k_B T \ln(K_d C_0)$$

How can we propagate the error δK_d into ΔG ?

$$\delta f = \left| \frac{\partial f}{\partial x} \right|_{x_0} \delta x$$

Substituting in the free energy definition above:

$$\delta \Delta G = \left| \frac{k_B T}{K_d C_0} \right| C_0 \delta K_d = k_B T \frac{\delta K_d}{K_d}$$



As we can see, a Gaussian with this standard deviation describes the uncertainty quite well!

Propagation of (t)error: Where does it come from, anyway?

If we have a function of multiple measurements or random variables, $f(x_1, x_2, \dots, x_N)$

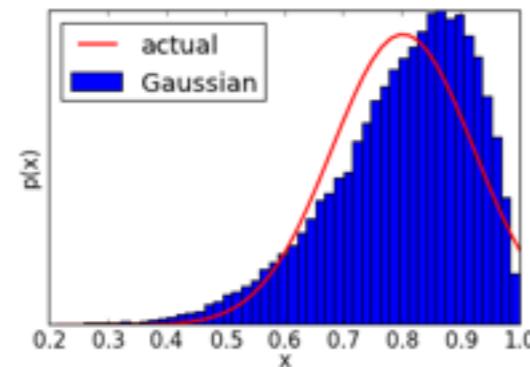
$$\delta^2 f = \sum_{i=1}^N \sum_{j=1}^N \frac{\partial f}{\partial x_i} \frac{\partial f}{\partial x_j} \Big|_{x_1, \dots, x_N} \delta x_i \delta x_j$$

If independent, $\delta x_i \delta x_j = 0$ for $i \neq j$

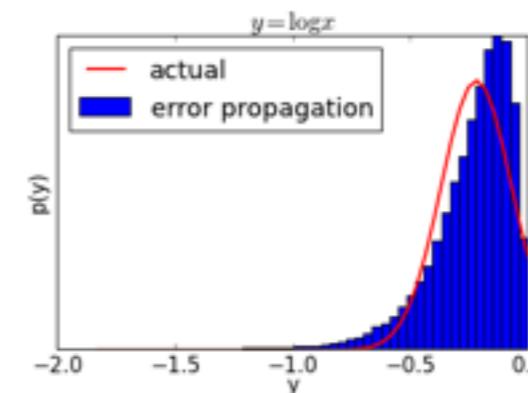
This can get totally messy as functions become complex or involve more variables.

CAUTION!

$$x \in (0, 1]$$



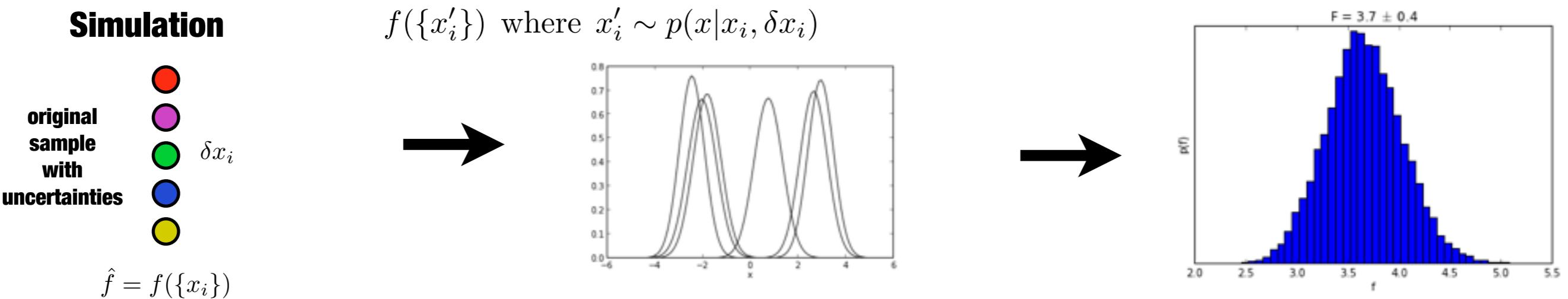
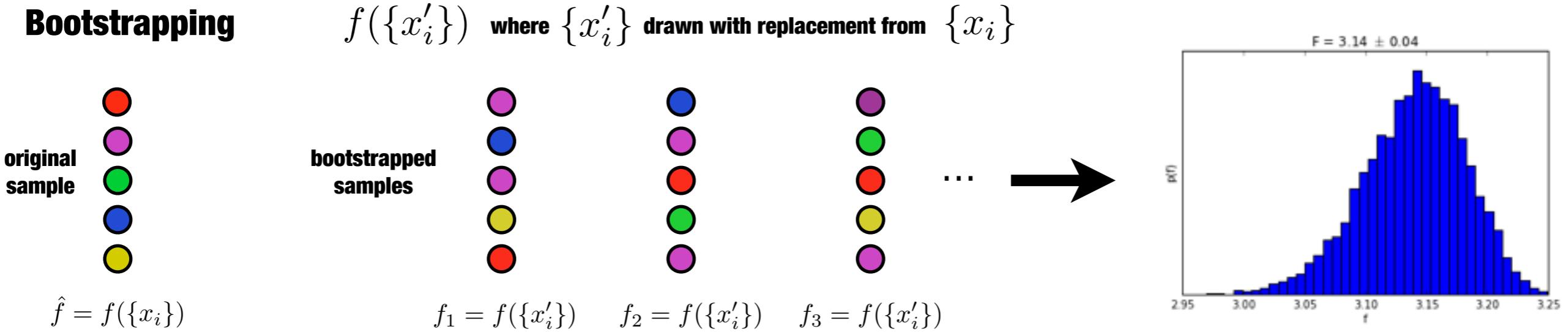
$$y = \log x$$



$$y \in (-\infty, 0]$$

For highly nonlinear functions, the first-order Taylor expansion doesn't do a good job.
Especially important for values on intervals, like probabilities ($p < 0$ and $p > 1$ nonsensical).

Three simple strategies for propagating error (without losing your mind)



Bayesian inference $\theta \sim p(\theta|\mathcal{D}) \propto p(\mathcal{D}|\theta) p(\theta)$

| | | | |
|------------|-----------|-----------------|-------|
| parameters | posterior | data likelihood | prior |
|------------|-----------|-----------------|-------|

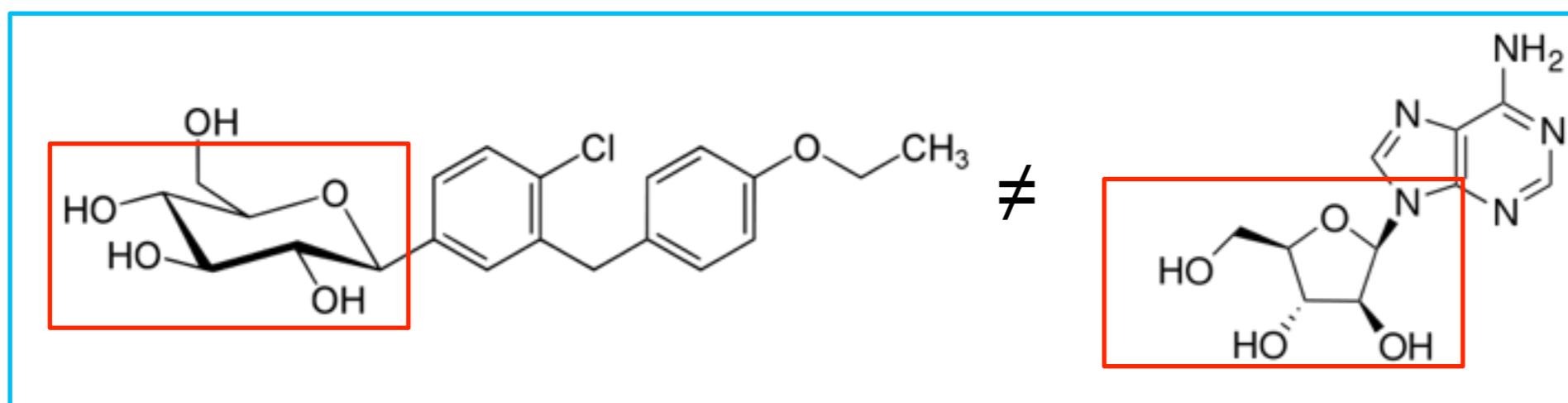
IPython tutorial available at <https://github.com/choderalab/cadd-grc-2013>

Noisy data: When is NO measurement or prediction actually a better option?

- When data will be over interpreted
 - When data is presented inappropriately
 - E.g. 23.427 ± 254.231
 - Example: Desperate chemists and solubility and PAMPA measurements
 - “I don’t care if the measurements aren’t quantitative, I need several significant figures”
 - “I don’t care if the answer is right, as long as it’s fast”
- When it will be used alone without related data
- When used out of context
- If not reported with sufficient caveats and error estimates

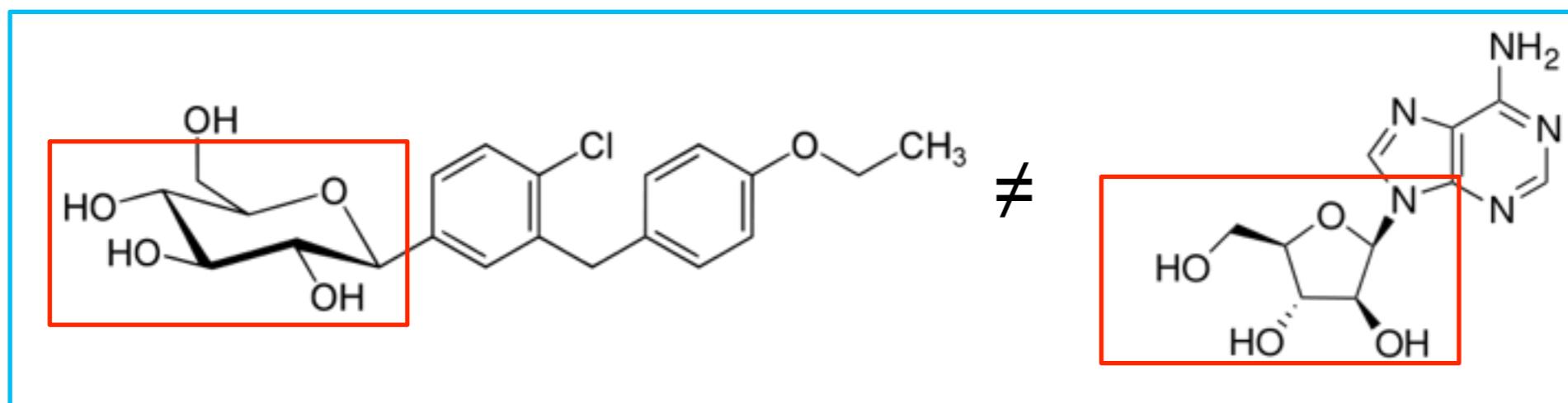
Noisy data: When is NO measurement or prediction actually a better option?

- Example: Over interpreted. without related data, out of context, without sufficient caveats and error estimates
- Blind over-interpretation of a mutagenicity alert from MultiCase toxicity prediction software that stalled a drug discovery program for 2 months



Noisy data: When is NO measurement or prediction actually a better option?

- Example: Over interpreted. without related data, out of context, without sufficient caveats and error estimates
- Blind over-interpretation of a mutagenicity alert from MultiCase toxicity prediction software that stalled a drug discovery program for 2 months



“There is no bad data, only bad data presentation”
“There are no bad soldiers, only bad generals”



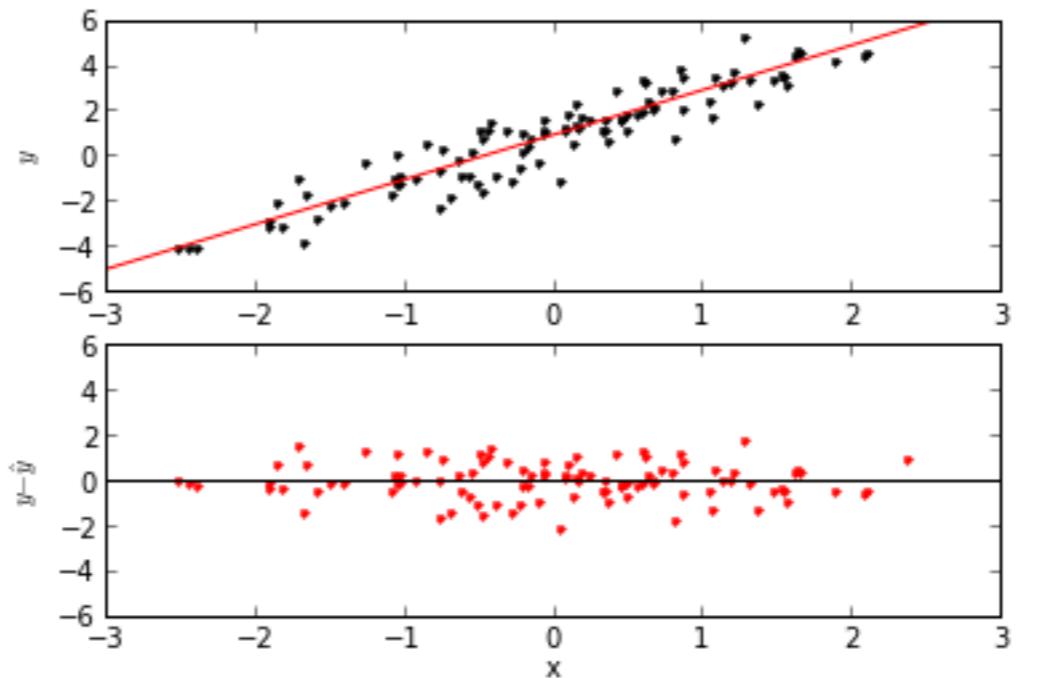
Noteworthy

- Models based on data of high confidence might be as good or better than many medium to high throughput assays
 - Given proper domain
 - Cyp 2C9
 - CACO-2
 - Met Stab
- Experimentalists could devote their time to high quality assays to provide data to enhance and improve models

Propagating error in simple models: A quick note on regression

$$\beta = \frac{\text{cov}(x, y)}{\text{var}(x)}$$

$$\alpha = \bar{y} - \beta \cdot \bar{x}$$



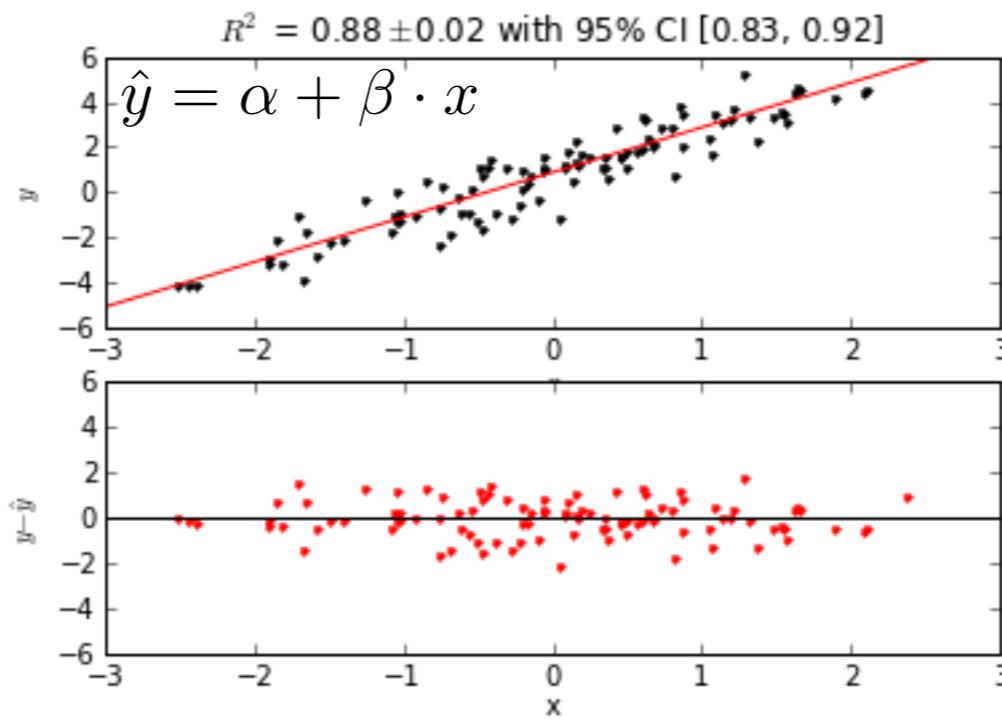
regression

residuals

Propagating error in simple models: Bootstrapping can estimate regression uncertainty

$$\beta = \frac{\text{cov}(x, y)}{\text{var}(x)}$$
$$\alpha = \bar{y} - \beta \cdot \bar{x}$$

$$R^2 = 1 - \frac{\hat{\sigma}_r^2}{\hat{\sigma}^2}$$



residual variance

$$\hat{\sigma}_r^2 = \frac{1}{N} \sum_{n=1}^N (y_i - \hat{y}_i)^2$$

sample variance

$$\hat{\sigma}_r^2 = \frac{1}{N} \sum_{n=1}^N (y_i - \bar{y})^2$$

Recall that uncertainty in coefficient of determination and fit parameters can be easily estimated by bootstrap.

Propagating error in simple models: Generalized R-squared

Cox & Snell pseudo-R²

$$R_{CS}^2 = 1 - \left(\frac{L(0)}{L(\hat{\theta})} \right)^{2/N}$$

residual variance

$$\hat{\sigma}_r^2 = \frac{1}{N} \sum_{n=1}^N (y_i - \hat{y}_i)^2$$

sample variance

$$\hat{\sigma}_r^2 = \frac{1}{N} \sum_{n=1}^N (y_i - \bar{y}_i)^2$$

Assume Gaussian error in deviations from model:

$$p(y|x) = (2\pi\sigma^2)^{-1/2} \exp \left[-\frac{(y - \hat{y}(x))^2}{2\sigma^2} \right]$$

Write likelihood for whole dataset:

$$L(\theta) = \prod_{i=1}^N (2\pi\sigma^2)^{-1/2} \exp \left[-\frac{(y_i - \hat{y}(x_i))^2}{2\sigma^2} \right] = (2\pi\sigma^2)^{-N/2} \exp \left[-\sum_{i=1}^N \frac{(y_i - \hat{y}(x_i))^2}{2\sigma^2} \right] = (2\pi\sigma^2)^{-N/2} e^{-N\hat{\sigma}_r^2/2\sigma^2}$$

If we don't know data error, we assume residual variance dominates:

$$\sigma^2 = \hat{\sigma}_r^2$$

This gives us likelihoods for fit model and null model (where data mean is assumed):

$$L(\hat{\theta}) = (2\pi\hat{\sigma}_r^2)^{-N/2} e^{-N/2}$$

$$L(0) = (2\pi\hat{\sigma}^2)^{-N/2} e^{-N/2}$$

Plugging these in, we recover the standard definition of the coefficient of determination!

$$R_{CS}^2 = 1 - \left(\frac{(2\pi\hat{\sigma}^2)^{-N/2} e^{-N/2}}{(2\pi\hat{\sigma}_r^2)^{-N/2} e^{-N/2}} \right)^{2/N} = \left(\frac{\hat{\sigma}_r^2}{\hat{\sigma}^2} \right) = R^2$$

Propagating error in simple models: Fitting data with differing errors

The regression problem can be generalized to the case where each measurement y_i has a corresponding uncertainty σ_i .

The likelihood is then

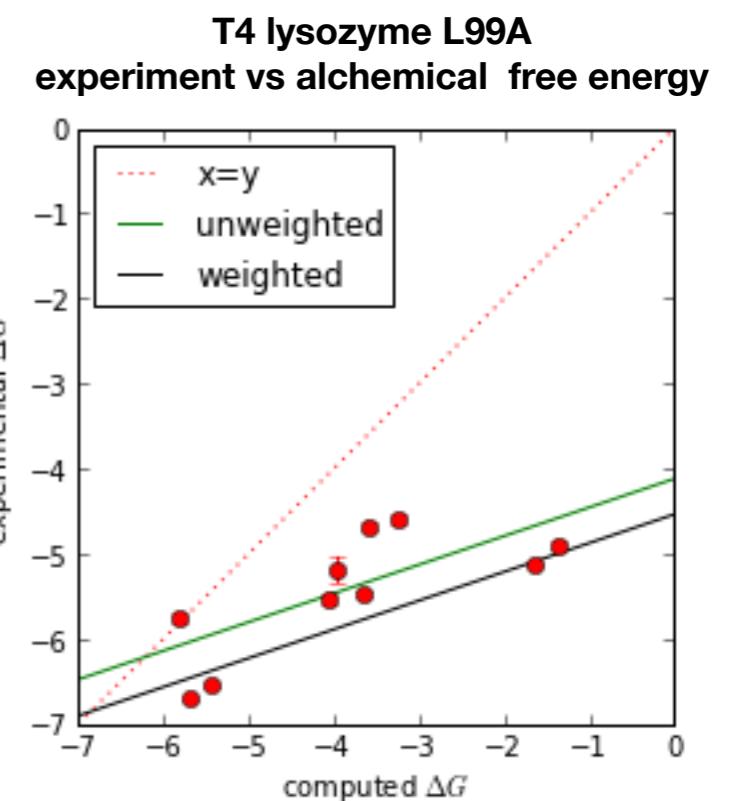
$$L(\theta) = \prod_{i=1}^N (2\pi\sigma_i^2)^{-1/2} \exp\left[-\frac{(y_i - \hat{y}(x_i))^2}{2\sigma_i^2}\right]$$

It is generally easier to maximize the *log-likelihood* instead:

$$l(\theta) = \sum_{i=1}^N \left[-\frac{1}{2} \log(2\pi\sigma_i^2) - \frac{(y_i - \hat{y}(x_i))^2}{2\sigma_i^2} \right] = -\frac{N}{2} \log(2\pi) + \sum_{i=1}^N \left[-\log \sigma_i^2 - \frac{(y_i - \hat{y}(x_i))^2}{2\sigma_i^2} \right]$$

To find the maximum likelihood estimate $\hat{\theta}$, we use the derivatives:

$$\frac{\partial}{\partial \theta} l(\theta) = \sum_{i=1}^N \frac{(y_i - \hat{y}(x_i))}{\sigma_i^2} \frac{\partial \hat{y}}{\partial \theta} \Big|_{x_i}$$



Outlier detection- when is it safe to drop data?

- Is there some reason other than it does not fit the current model?
 - Suspicious or simply inconvenient?
- Is the data “bad”?
 - Is it warranted by the experimental metadata?
 - Age of compound / Date of assay (vs. that of other compounds) / Purity / Assay conditions / Source
 - Are the physical properties suspicious (calculated or actual)
 - Reactivity / Promiscuity / Solubility
- Is the measurement or compound just different?
 - The compound
 - logP, MW, # rotatable bonds, amount of stereo centers, etc
 - The measurement
 - Leverage – is the value way outside of bounds?
 - Assay variation

What can we do about outliers?

Noise accommodation vs. outlier detection

Outliers are data that are **highly improbable** for the model under consideration. They often represent real-world anomalies, and may be interesting or useless.

Outliers may be due to:

- change in environment/assay conditions
- unappreciated physical phenomena
- instrumentation error
- human error (mishandling, data entry errors)
- incorrect conversions
- numeric codes to represent missing data

We deal with outliers in one of two ways:

1. If our model can't handle outliers, we use **noise removal** or **noise accommodation** [e.g. robust statistics, robust regression] c.f. Maria Gallardo's dipole moments
2. If the outliers are important phenomena we're trying to identify: **outlier detection** [e.g. uncharacterized effects, significant events, compounds to retest]

Sometimes, frequentists can be handy

A simple statistical tests for outliers: Grubbs' test

Assume **normality of data distribution; sequentially identify outliers.**

Grubbs' test statistic: $G = \frac{\max_{i=1,\dots,N} |y_i - \bar{y}|}{s}$

| | |
|-----------|--------------------------------|
| s | sample standard deviation |
| \bar{y} | sample mean |
| α | significance level (e.g. 0.05) |

There are outliers if we find $G > G_{\text{test}}$

$$G_{\text{test}} = \frac{N-1}{\sqrt{N}} \sqrt{\frac{t^2}{N-2+t^2}}$$

t is the upper critical value of t-distribution with N-2 dof and significance level $\alpha/(2N)$
(Bonferroni correction)

```
from scipy import stats
alpha = 0.95 # 95% significance level
t = stats.t.isf(1-alpha/(2*N), N-2)
Gtest = (N-1)/sqrt(N)*sqrt(t**2 / (N-2+t**2))
```

Statistical tests for outliers: Grubbs' test

An example from your nightmares



 Technical Note

Does Your Lab Coat Fit to Your Assay?



Michael Busch¹, Heinz Bjoern Thoma¹, and Ingo Kober¹

Journal of Biomolecular Screening
18(6) 744–747
© 2013 Society for Laboratory
Automation and Screening
DOI: 10.1177/1087057113481621
jbx.sagepub.com



Abstract

An explanation for randomly occurring spikes on microplates in fluorescence-based assays employing shorter-wavelength readouts is presented. It is demonstrated that lint originating from standard (white cotton) lab coats is most likely to be responsible for such artifacts in assays applying wavelengths at 380 nm excitation and 450 nm emission. The fluorescence properties of this lint are discussed and compared with those of optical brighteners. An alternative to the use of cotton-based lab coats is presented, which led to a reduction of spikes in a high-throughput screening campaign by 90%.

Keywords

HTS, fluorescence, spikes, lab coat, dust, lint

“During assay development...about 5 to 10 randomly distributed spikes (signals up to 8 times higher than a normal assay signal) were found on each 384-well plate. This would have led to 1% to 2% of false-positives in HTS...”

Statistical tests for outliers: Grubbs' test

An example from your nightmares



empty microplate fresh out of packaging

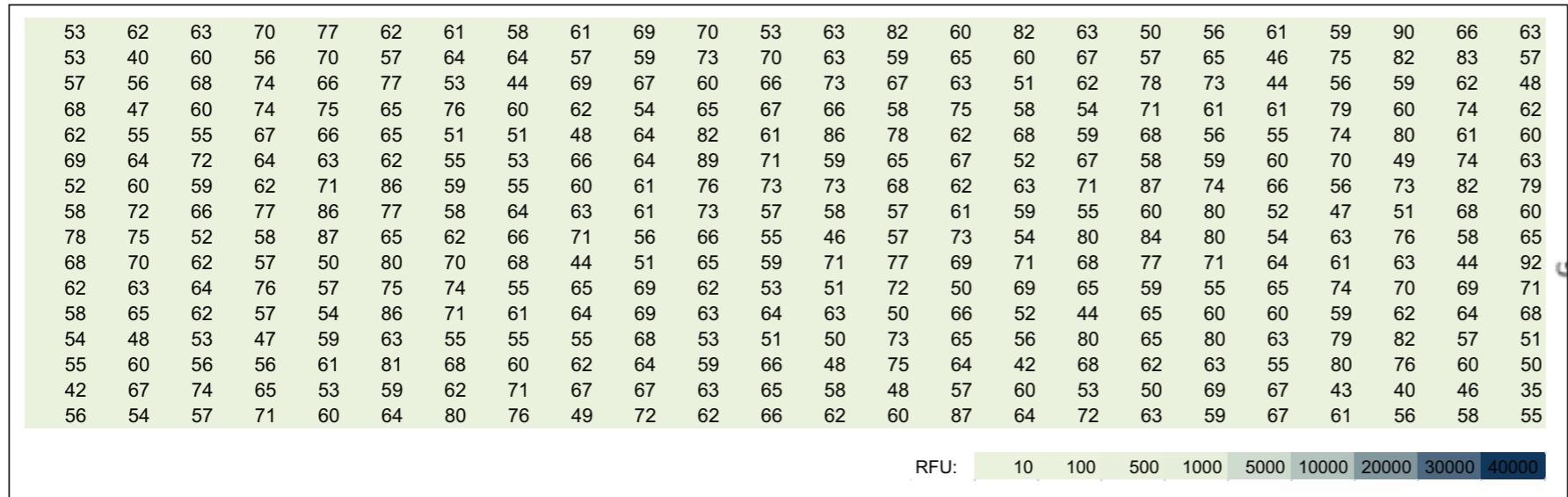
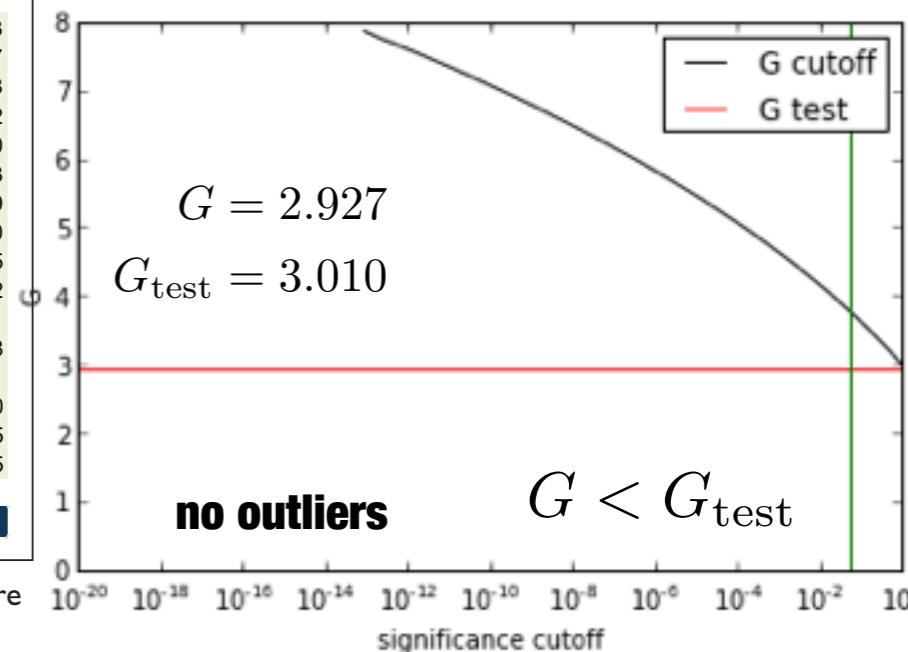


Figure 1. Empty microplate, freshly taken out of a closed bag. Fluorescence readout at 380/450 nm ex/em. The medium signals were less than 70 RFU.



empty microplate exposed to laboratory air for 2 min

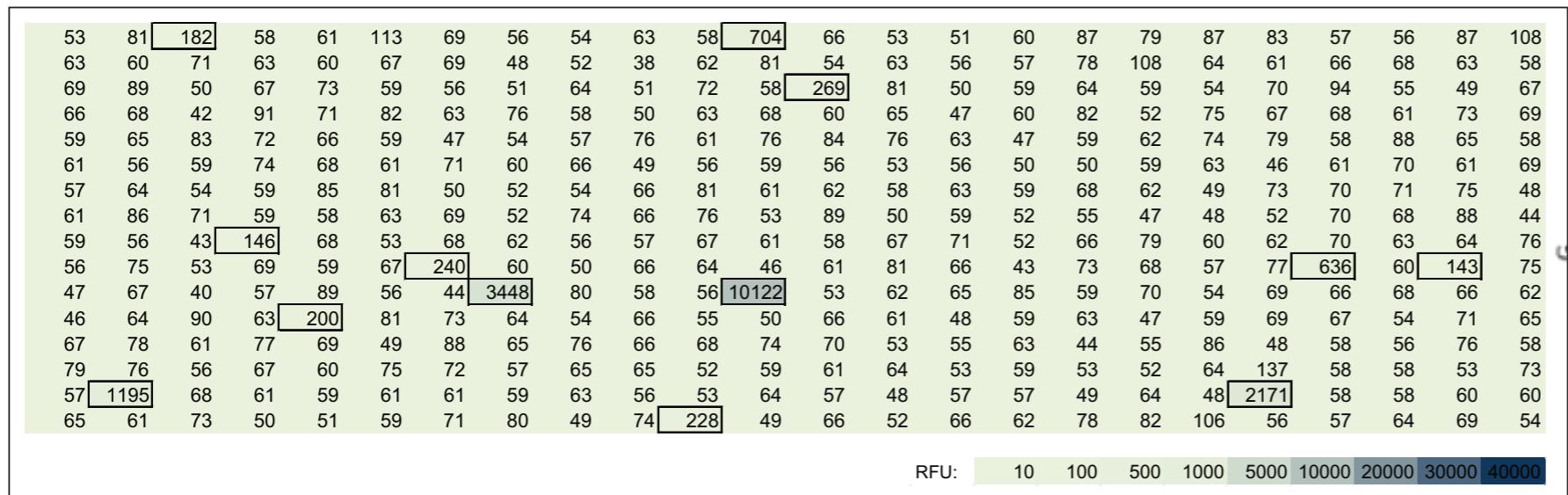
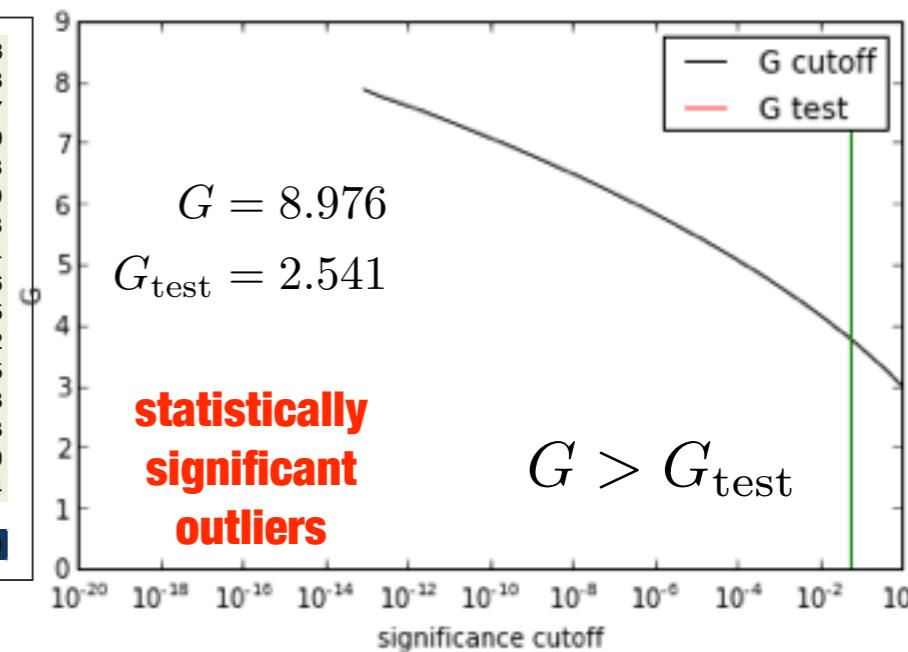


Figure 2. Measurement of the same plate after moving it for 2 min in free laboratory air. Spikes up to 10 000 RFU were observed here.



Statistical tests for outliers: Grubbs' test

An example from your nightmares

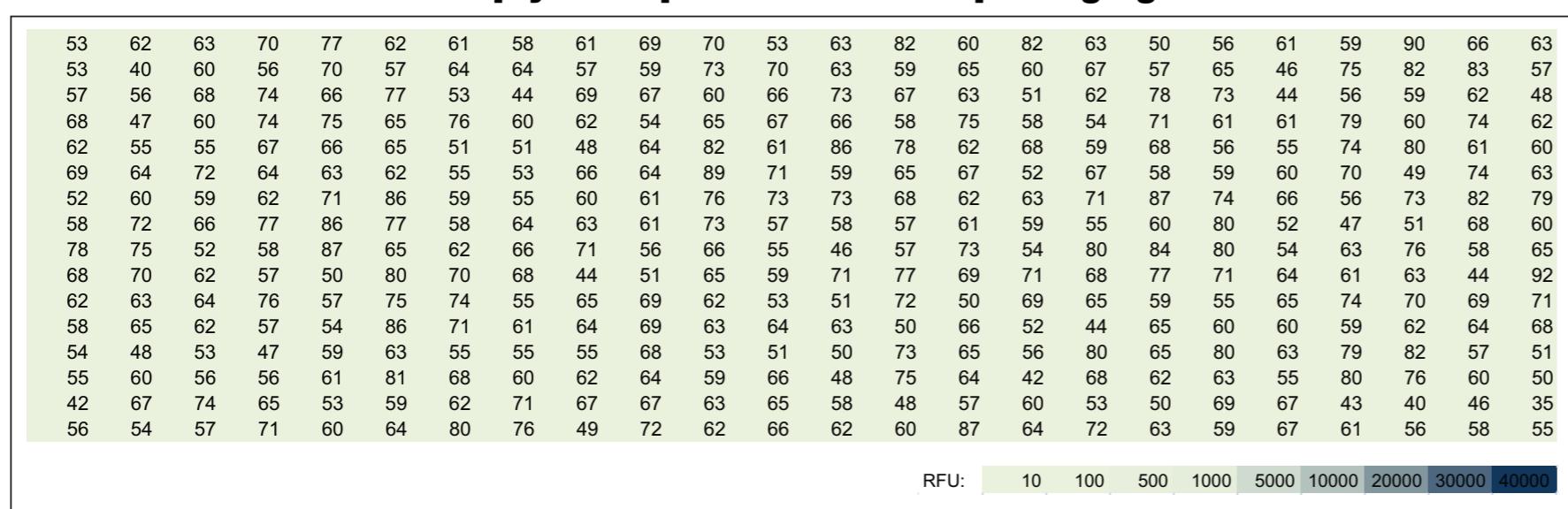


Figure 1. Empty microplate, freshly taken out of a closed bag. Fluorescence readout at 380/450 nm ex/em. The medium signals were less than 70 RFU.

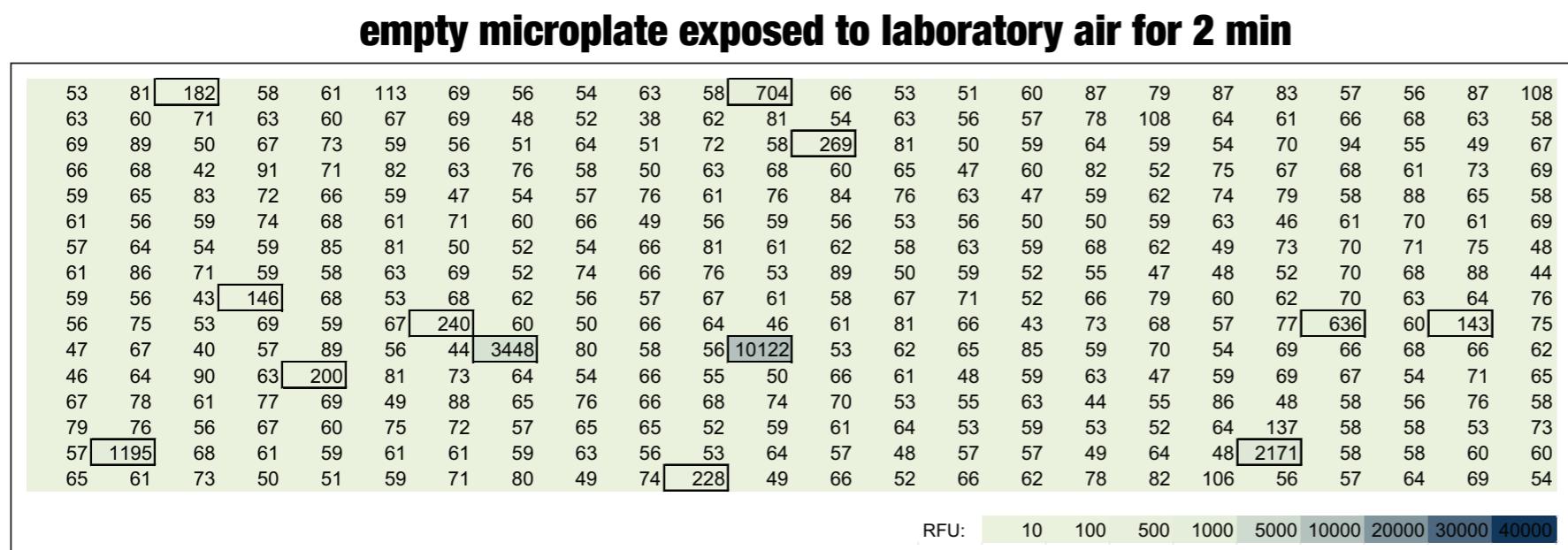
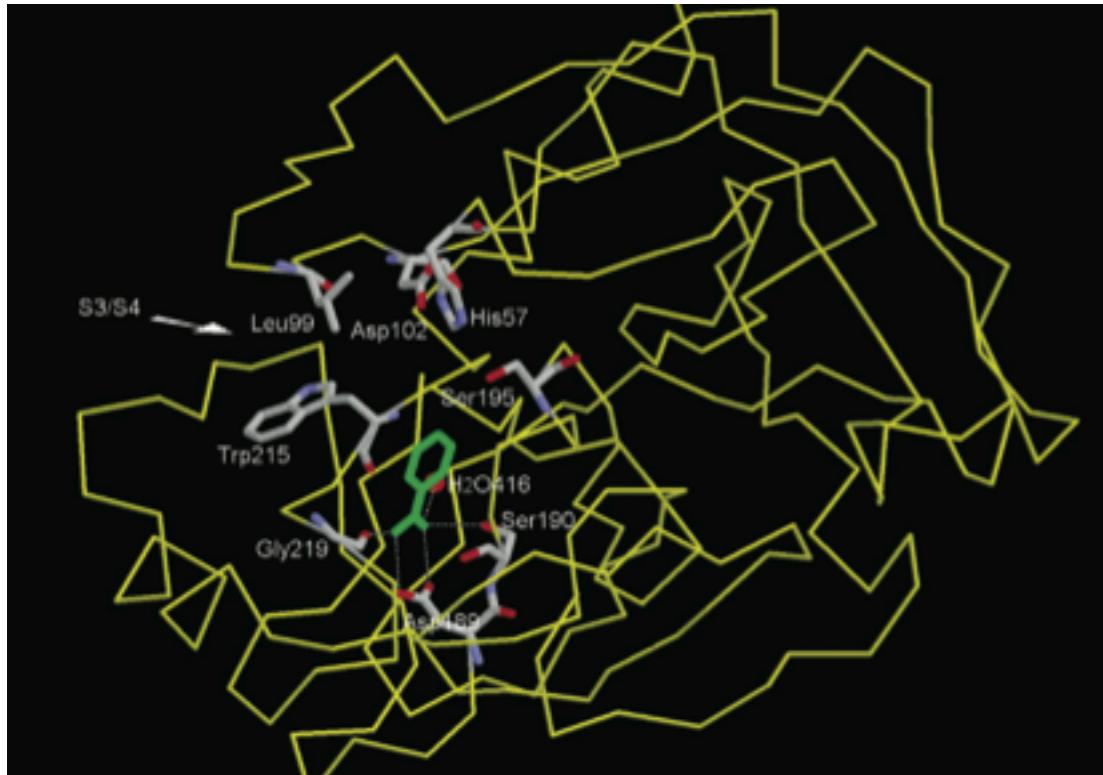


Figure 2. Measurement of the same plate after moving it for 2 min in free laboratory air. Spikes up to 10 000 RFU were observed here.



Dark tales of woe: Entropy-enthalpy compensation: Fact or fiction?

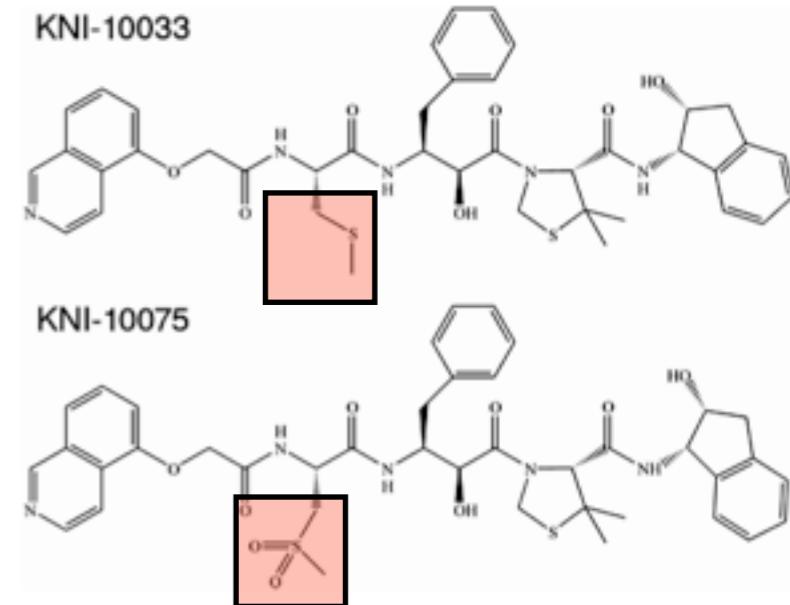
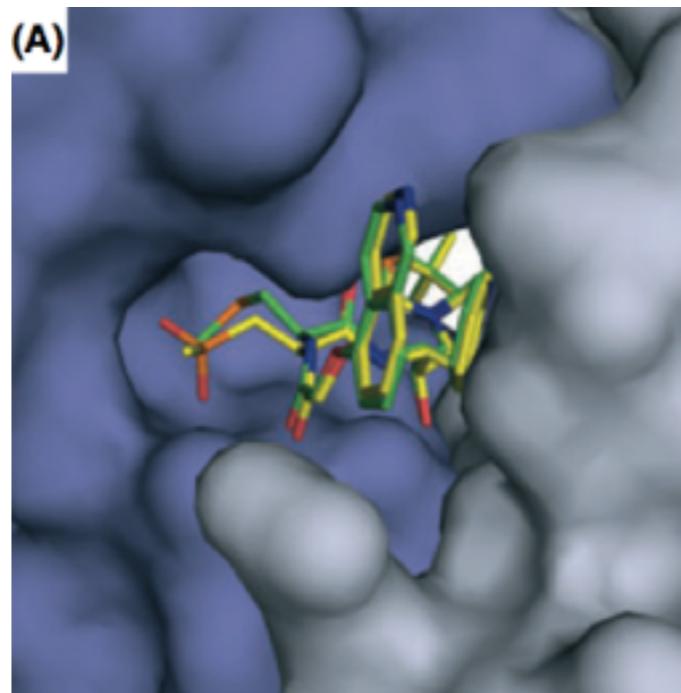
p-alkylbenzamidinium inhibitors of trypsin



| | ΔG (kJ/mol) | ΔH | $T\Delta S$ |
|----------|---------------------|-------------|-------------|
| R = H | -26.6 ± 0.3 | -18.9 ± 0.4 | 7.7 ± 0.6 |
| R = n-Bu | -26.2 ± 0.3 | -9.9 ± 0.5 | 16.3 ± 0.7 |

Small chemical modifications appear to drastically change entropic and enthalpic components, but have little effect on affinity.

HIV-1 protease inhibitors

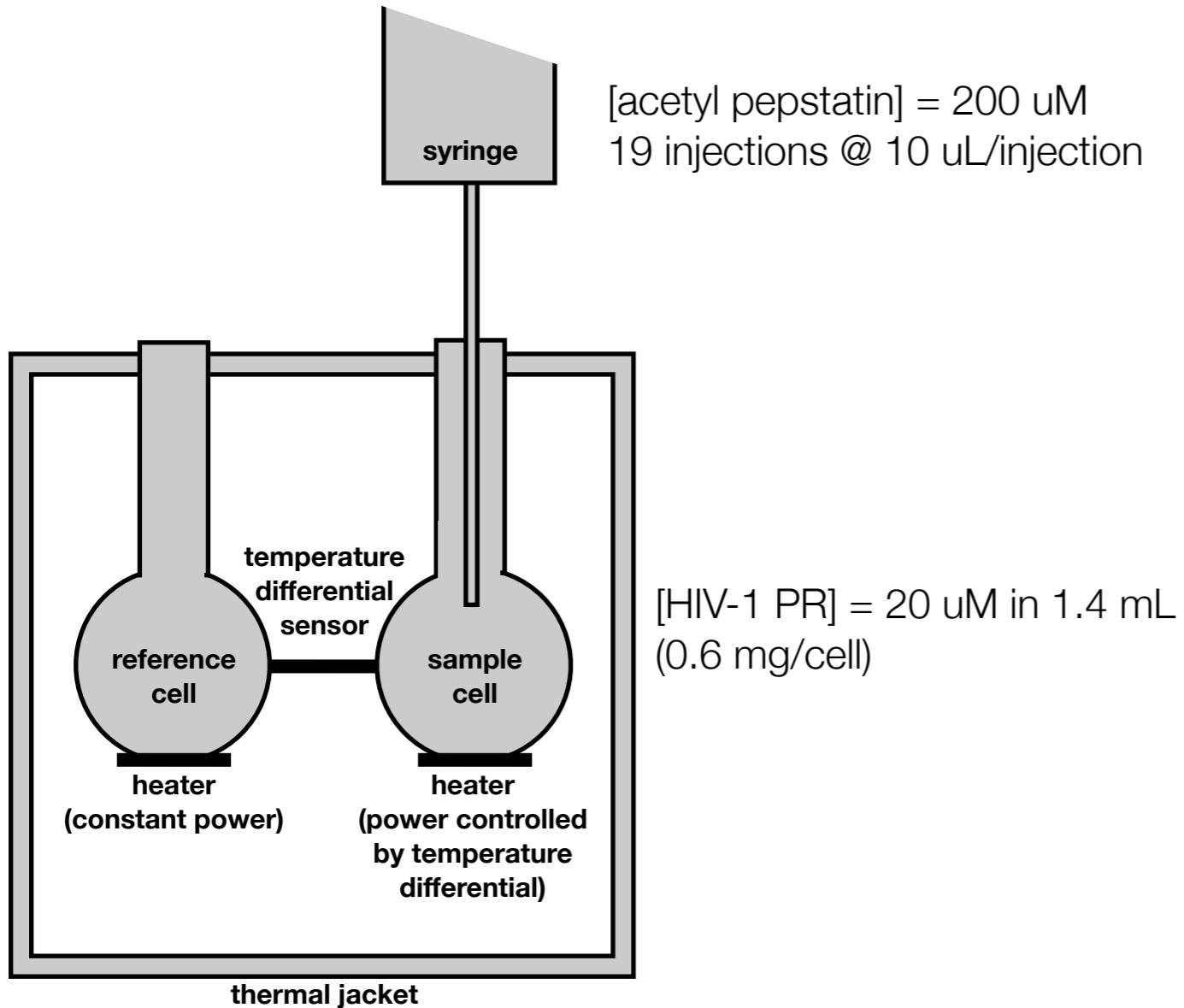


| | ΔG (cal/mol) | ΔH (cal/mol) | $-T\Delta S$ (cal/mol) | K_d (M) |
|-----------------|----------------------|----------------------|------------------------|-----------------------------|
| KNI-10033 → pWT | -14 870 ± 90 | -8200 ± 230 | -6670 ± 90 | $1.3^{-11} \pm 2 10^{-12}$ |
| KNI-10075 → pWT | -14 620 ± 190 | -12 120 ± 610 | -2500 ± 190 | $2 10^{-11} \pm 8 10^{-12}$ |

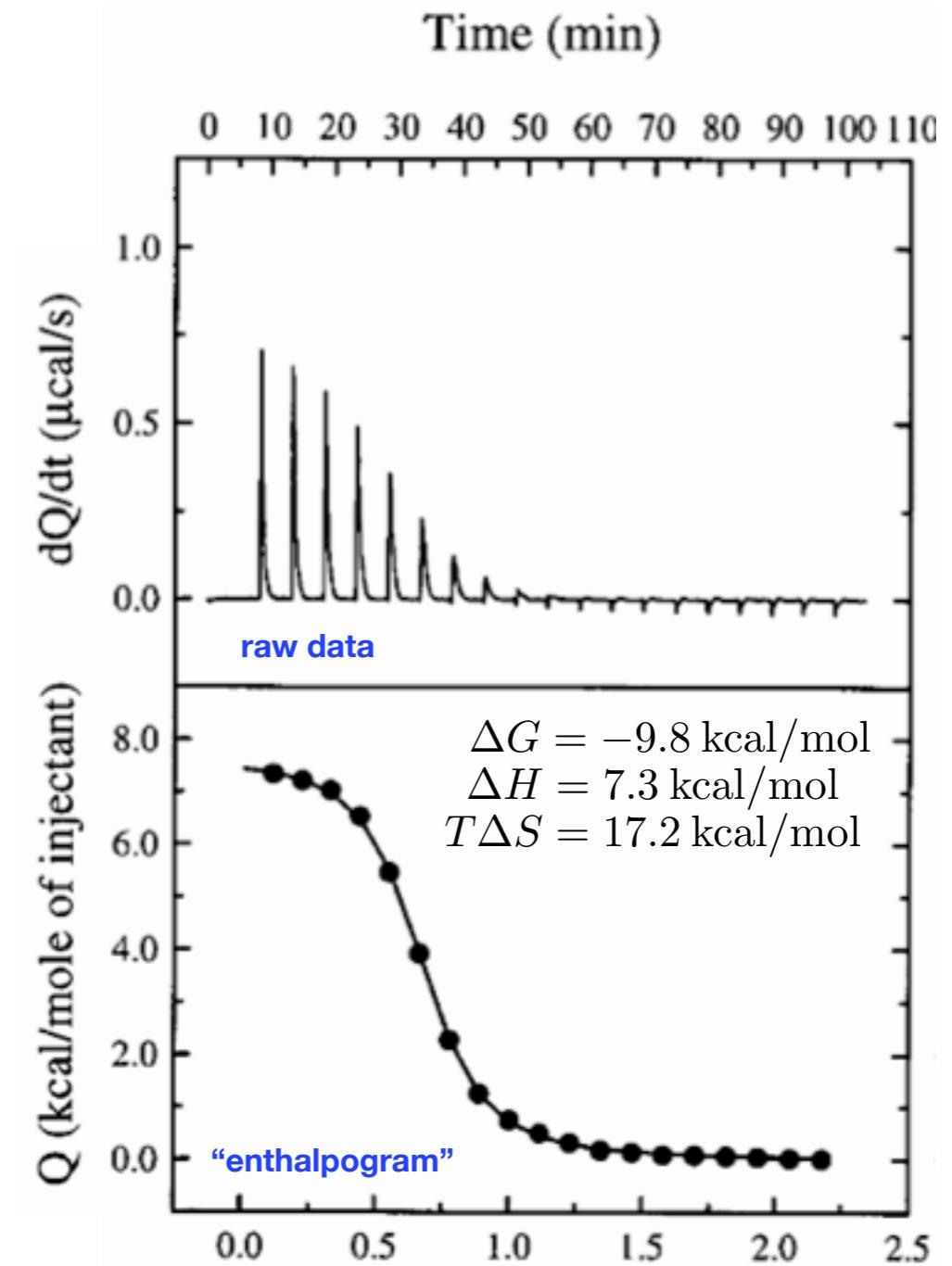
The enthalpy changes were measured at 15, 25 and 35 °C to determine the heat capacity change. (ΔC_p) is equal to -461 ± 34 cal/K/mol for KNI-10033 and -305 ± 41 cal/K/mol for KNI-10075.

Isothermal titration calorimetry (ITC) can simultaneously interrogate free energies and enthalpies of binding

[illustrative experiment]



[HIV-1 PR] = 20 uM in 1.4 mL
(0.6 mg/cell)



(Note that some reactions have no measurable change in heat, and are not measurable by ITC.)

Velazquez-Campoy A, Kiso Y, and Friere E. Arch. Biochem. Biophys. 390:169, 2001.

How reliable is calorimetric data, really?

A test of variation among truly independent experiments

The ABRF-MIRG'02 study:

Send identical aliquots of the **same sample** of protein and ligand to 14 core analysis facilities (experts!) and ask them to report the measured ΔG and ΔH by ITC.

The should get the **same answer**, within error.

The reported errors should match the variation among the reported results.

This experiment is almost never repeated because of the large quantity of protein needed for one ITC experiment, and the undesirability of **repeating the experiment from scratch** multiple times.

This is pretty much the only dataset of its kind reported in the literature.

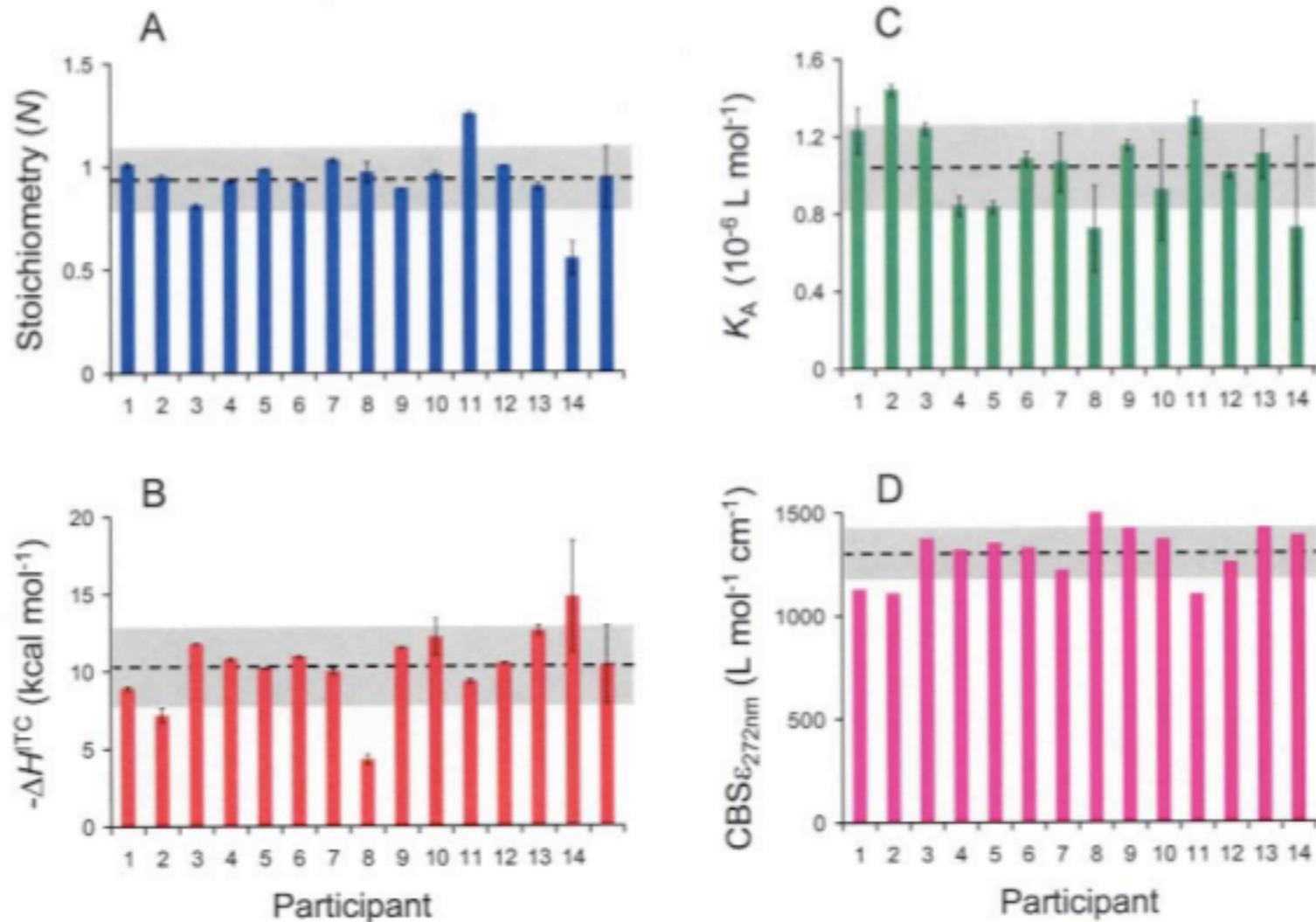


FIGURE 5

ITC characterization of the CBS/CA II interaction (based on Fig. 4 and Table 3). **A:** Stoichiometry (0.94 ± 0.15). **B:** Enthalpy upon binding ($-10.4 \pm 2.5 \text{ kcal mol}^{-1}$). **C:** Affinity [$(1.00 \pm 0.22) \times 10^6 \text{ L mol}^{-1}$]. **D:** molar extinction coefficient for CBS at 272 nm ($1307 \pm 126 \text{ L mol}^{-1} \text{ cm}^{-1}$). Mean values and the standard deviation for 14 determinations are represented by the horizontal dotted lines and gray bands. Error bars denote the standard deviation of nonlinear least squares analysis, except for participants 10 and 14 who reported standard deviations for replicate analyses. No errors were reported for the extinction coefficient determinations.

How reliable is calorimetric data, really? Error is significantly underreported

TABLE 3

Summary of Reported Isothermal Titration Calorimetry Results^a

| Participant | Molar Binding Ratio (<i>N</i>) | K_A ($10^{-6} \times L\ mol^{-1}$) | ΔH^{ITC} (kcal mol $^{-1}$) | C-value ^b | Control Titrations |
|-----------------|----------------------------------|--|---|----------------------|---|
| 1 | 1.01 ± 0.01 | 1.2 ± 0.1 | -8.9 ± 0.1 | 17 | CBS into buffer and buffer into CA II. Former subtracted as part of data analysis. |
| 2 | 0.95 ± 0.01 | 1.44 ± 0.03 | -7.2 ± 0.5 | 75 | None |
| 3 | 0.81 ± 0.01 | 1.24 ± 0.03 | -11.8 ± 0.02 | 44 | CBS into buffer and buffer into CA II. Former subtracted as part of data analysis. |
| 4 | 0.929 ± 0.007 | 0.84 ± 0.05 | -10.8 ± 0.1 | 24 | Pilot run. |
| 5 | 0.987 ± 0.003 | 0.84 ± 0.03 | -10.20 ± 0.04 | 41 | CBS into buffer, subtracted in data analysis. |
| 6 | 0.921 ± 0.003 | 1.08 ± 0.04 | -10.95 ± 0.05 | 39 | CBS into buffer, subtracted in data analysis. |
| 7 | 1.03 ± 0.01 | 1.1 ± 0.2 | -10.0 ± 0.2 | 55 | CBS into buffer, subtracted in data analysis. |
| 8 | 0.97 ± 0.05 | 0.7 ± 0.2 | -4.3 ± 0.3 | 7.0 | CBS into buffer and buffer into CA II. Both subtracted as part of data analysis. |
| 9 | 0.891 ± 0.002 | 1.15 ± 0.03 | -11.53 ± 0.04 | 59 | CBS into buffer, average value subtracted in analysis. |
| 10 ^c | 0.96 ± 0.02 | 0.9 ± 0.2 | -12 ± 1 | 34 | CBS into buffer and buffer into CA II. Average of former used in analysis. |
| 11 | 1.25 ± 0.01 | 1.3 ± 0.1 | -9.3 ± 0.1 | 9.0 | None |
| 12 | 1.000 ± 0.003 | 1.01 ± 0.03 | -10.51 ± 0.04 | 29 | Buffer into buffer. |
| 13 | 0.90 ± 0.02 | 1.1 ± 0.1 | -12.6 ± 0.3 | 12 | CBS into buffer, linear regression subtracted in analysis. |
| 14 ^c | 0.55 ± 0.08 | 0.7 ± 0.3 | -15 ± 4 | 22 | CBS into buffer, subtracted in analysis. |

Note that observed 20% error in K_A gives only ± 0.1 kcal/mol error in ΔG , while 20% error in ΔH directly impacts ΔH .

This means absolute error in ΔG is actually still small, while absolute error in ΔH is big.

**The reported error bars cannot be trusted.
They're often an order of magnitude (or more!) too small.**

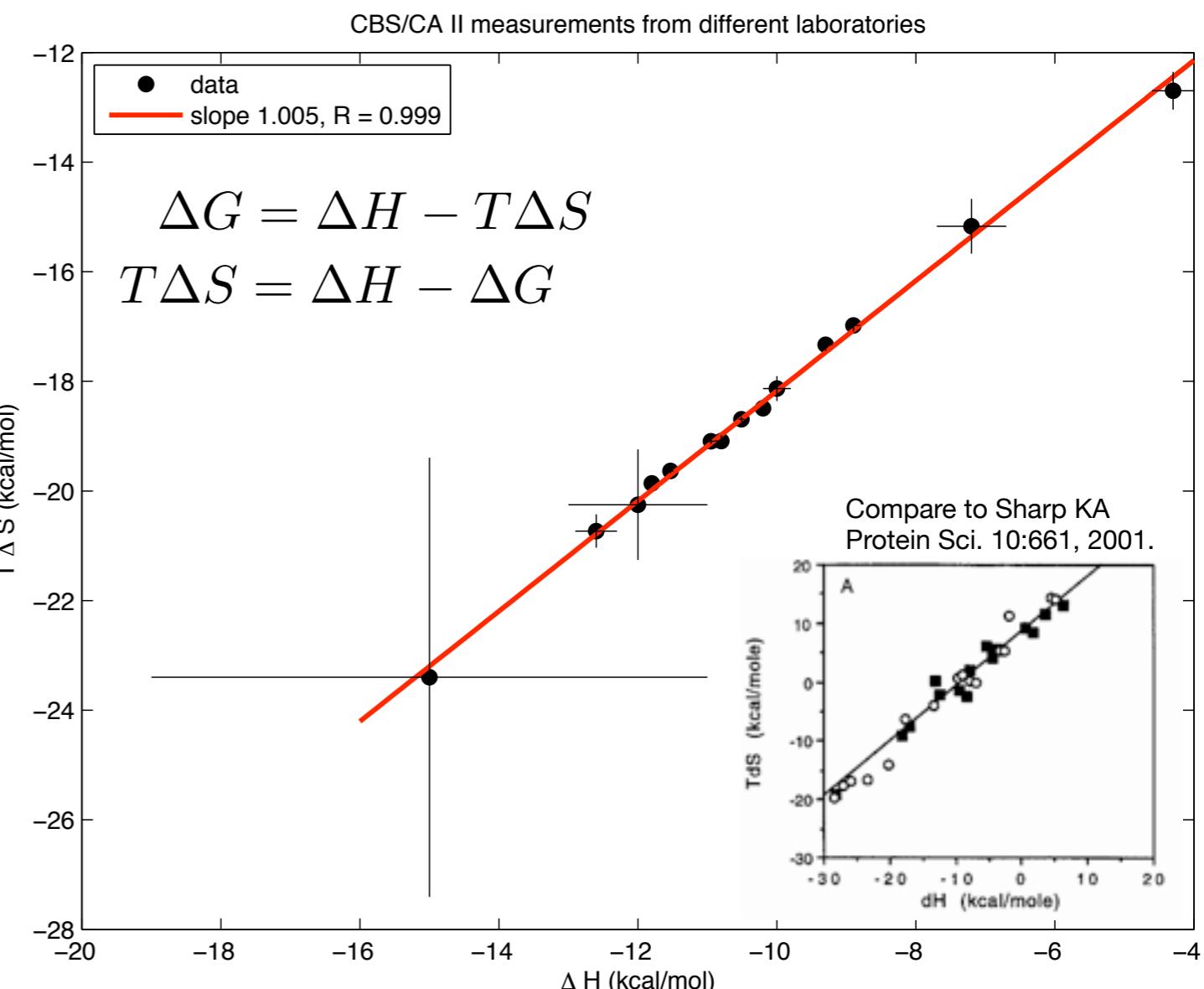
ITC measurements on the **same** system by different laboratories exhibits apparent entropy-enthalpy compensation!

TABLE 3

Summary of Reported Isothermal Titration Calorimetry Results^a

| Participant | Molar Binding Ratio (<i>N</i>) | K_A ($10^{-4} \times L \text{ mol}^{-1}$) | ΔH^{ITC} (kcal mol $^{-1}$) | C-value ^b | Control Titrations |
|-----------------|----------------------------------|---|---|----------------------|--|
| 1 | 1.01 ± 0.01 | 1.2 ± 0.1 | -8.9 ± 0.1 | 17 | CBS into buffer and buffer into CA II. Former subtracted as part of data analysis. |
| 2 | 0.95 ± 0.01 | 1.44 ± 0.03 | -7.2 ± 0.5 | 75 | None |
| 3 | 0.81 ± 0.01 | 1.24 ± 0.03 | -11.8 ± 0.02 | 44 | CBS into buffer and buffer into CA II. Former subtracted as part of data analysis. |
| 4 | 0.929 ± 0.007 | 0.84 ± 0.05 | -10.8 ± 0.1 | 24 | Pilot run. |
| 5 | 0.987 ± 0.003 | 0.84 ± 0.03 | -10.20 ± 0.04 | 41 | CBS into buffer, subtracted in analysis. |
| 6 | 0.921 ± 0.003 | 1.08 ± 0.04 | -10.95 ± 0.05 | 39 | CBS into buffer, subtracted in analysis. |
| 7 | 1.03 ± 0.01 | 1.1 ± 0.2 | -10.0 ± 0.2 | 55 | CBS into buffer, subtracted in analysis. |
| 8 | 0.97 ± 0.05 | 0.7 ± 0.2 | -4.3 ± 0.3 | 7.0 | CBS into buffer, subtracted in analysis. Both subtracted as part of data analysis. |
| 9 | 0.891 ± 0.002 | 1.15 ± 0.03 | -11.53 ± 0.04 | 59 | CBS into buffer, average value subtracted in analysis. |
| 10 ^c | 0.96 ± 0.02 | 0.9 ± 0.2 | -12 ± 1 | 34 | CBS into buffer and buffer into CA II. Average of former used in analysis. |
| 11 | 1.25 ± 0.01 | 1.3 ± 0.1 | -9.3 ± 0.1 | 9.0 | None |
| 12 | 1.000 ± 0.003 | 1.01 ± 0.03 | -10.51 ± 0.04 | 29 | Buffer into buffer. |
| 13 | 0.90 ± 0.02 | 1.1 ± 0.1 | -12.6 ± 0.3 | 12 | CBS into buffer, linear regression subtracted in analysis. |
| 14 ^c | 0.55 ± 0.08 | 0.7 ± 0.3 | -15 ± 4 | 22 | CBS into buffer, subtracted in analysis. |

plot



Myszka et al. J. Biomol. Tech. 14:247, 2003.

How can these plots be real evidence of compensation if we can generate the same plot from different measurements on same system?

Analysis of ITC experiments: The Bayesian way

$$p(\theta|\mathcal{D}) \propto p(\mathcal{D}|\theta)p(\theta)$$

\mathcal{D} data

θ model parameters

$p(\theta|\mathcal{D})$ posterior

$p(\mathcal{D}|\theta)$ sampling distribution (model)

$p(\theta)$ prior

Analysis of ITC experiments: The Bayesian way

$$p(\theta|\mathcal{D}) \propto p(\mathcal{D}|\theta)p(\theta)$$

\mathcal{D} data

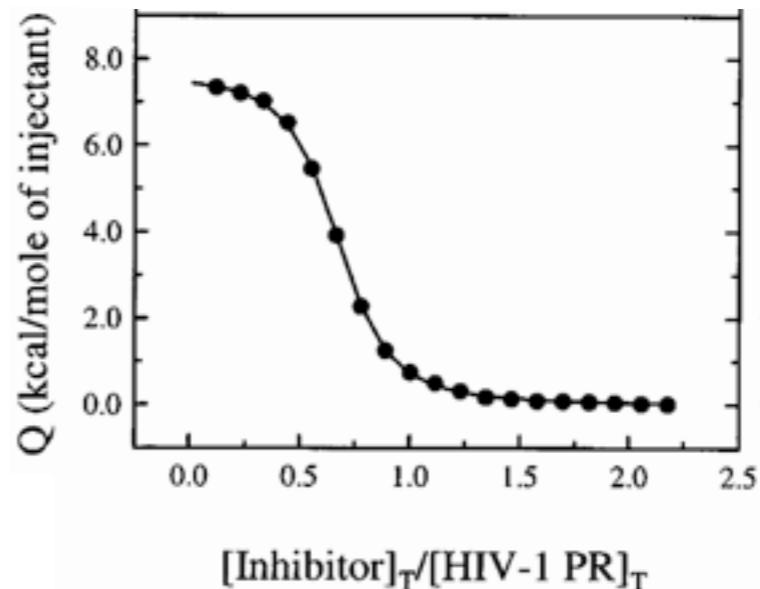
$\mathcal{D} = \{q_1, q_2, \dots, q_N\}$ measurements of evolved heat

θ model parameters

$p(\theta|\mathcal{D})$ posterior

$p(\mathcal{D}|\theta)$ sampling distribution (model)

$p(\theta)$ prior



Analysis of ITC experiments: The Bayesian way

$$p(\theta|\mathcal{D}) \propto p(\mathcal{D}|\theta)p(\theta)$$

| | | | |
|-------------------------|-------------------------------|--|-----------------------------------|
| \mathcal{D} | data | $\theta = \{\Delta G, \Delta H, T\Delta S, \Delta H_0\}$ | thermodynamic parameters |
| θ | model parameters | | |
| $p(\theta \mathcal{D})$ | posterior | $\Delta G = -kT \ln K_a$ | free energy of binding |
| $p(\mathcal{D} \theta)$ | sampling distribution (model) | ΔH | enthalpy of binding |
| $p(\theta)$ | prior | $T\Delta S$ | entropic contribution to binding |
| | | ΔH_0 | heat of dilution |
| | | | ... any additional parameters ... |

Analysis of ITC experiments: The Bayesian way

$$p(\theta|\mathcal{D}) \propto p(\mathcal{D}|\theta)p(\theta)$$

| | | | |
|-------------------------|-------------------------------|---|---|
| \mathcal{D} | data | $p(\mathcal{D} \theta) = \prod_{n=1}^N \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{(q_n - q_n^*)^2}{2\sigma^2}\right]$ | Gaussian error model |
| θ | model parameters | | |
| $p(\theta \mathcal{D})$ | posterior | q_n | measured heat of injection n |
| $p(\mathcal{D} \theta)$ | sampling distribution (model) | q_n^* | true heat of injection n |
| $p(\theta)$ | prior | σ | std dev of error in measured heat (nuisance parameter) |

$$P + L \xrightarrow{\Delta H} PL$$

$$q_n^* = Q_n - Q_{n-1}$$

$$Q_n = \Delta H \cdot V_n [PL]_n + n\Delta H_0 \quad \text{heat potential}$$

Analysis of ITC experiments: The Bayesian way

$$p(\theta|\mathcal{D}) \propto p(\mathcal{D}|\theta)p(\theta)$$

| | | | |
|-------------------------|-------------------------------|---|---|
| \mathcal{D} | data | $p(\Delta G, \Delta H, \Delta H_0, \sigma) \propto \sigma^{-1}$ | prior on measurement noise |
| θ | model parameters | $[M]_0^* \sim \mathcal{N}([M]_0, \sigma_M^2)$ | prior on cell and syringe concentrations |
| $p(\theta \mathcal{D})$ | posterior | $[L]_s^* \sim \mathcal{N}([L]_s, \sigma_L^2)$ | |
| $p(\mathcal{D} \theta)$ | sampling distribution (model) | | |
| $p(\theta)$ | prior | $\Delta G, \Delta H, \Delta H_0$ | can be of any sign and value |
| | | $\sigma > 0$ | scale parameter; can be of any magnitude (Later, could build in some <i>a priori</i> knowledge of instrument error or calibration runs.) |

Important points:

- * We don't know the size of the measurement error
 - * We don't know the actual ligand and protein concentrations
- No problem: Just make them **nuisance parameters!**

Analysis of ITC experiments: The Bayesian way

$$p(\theta|\mathcal{D}) \propto p(\mathcal{D}|\theta)p(\theta)$$

$$p(\theta|\mathcal{D}) = (2\pi)^{-N/2}\sigma^{-(N+1)} \exp\left[-\frac{1}{2\sigma^2} \sum_{n=1}^N (q_n - q_n^*)^2\right] \text{ posterior}$$

\mathcal{D} data
 θ model parameters

$p(\theta|\mathcal{D})$ posterior

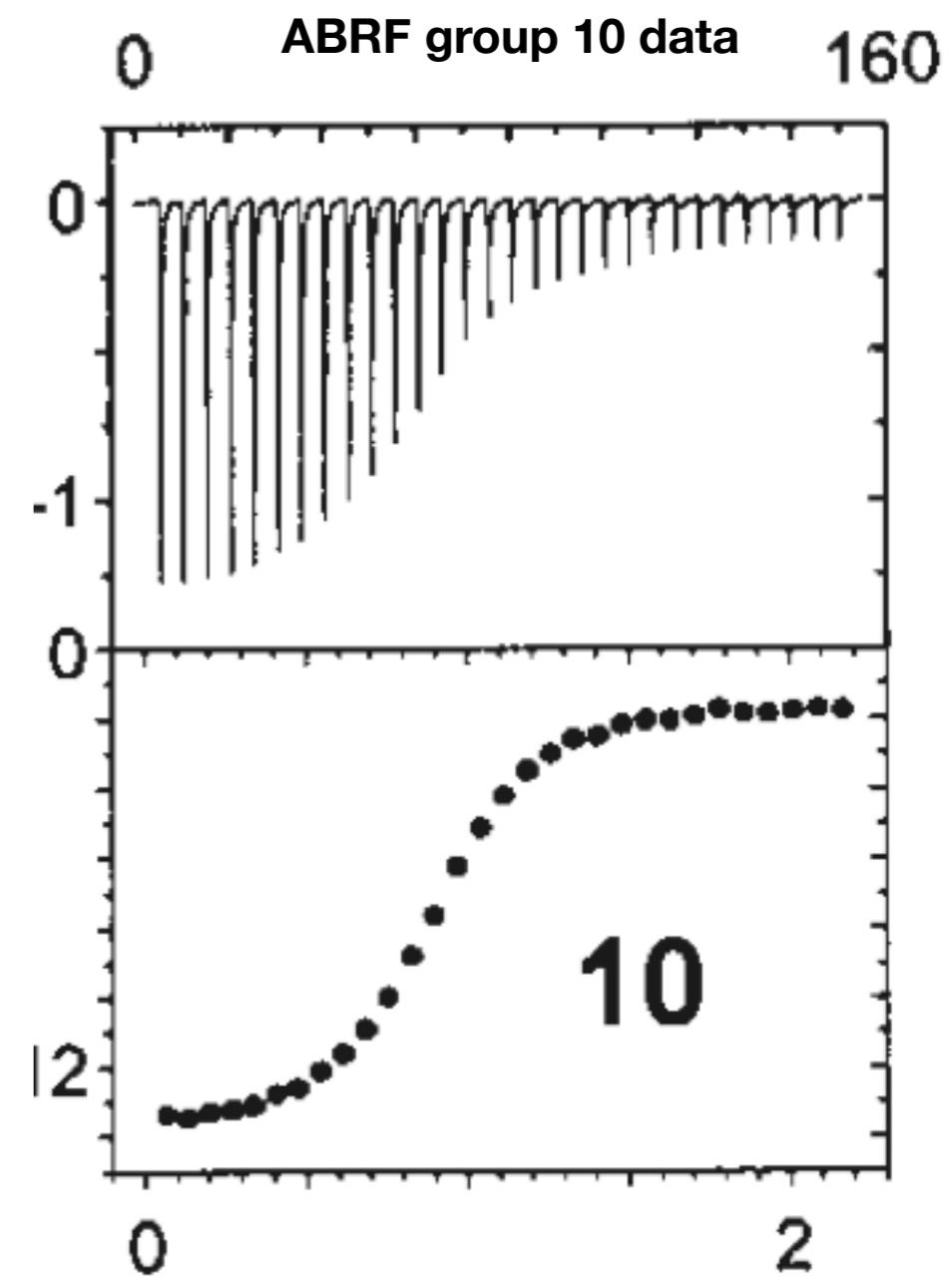
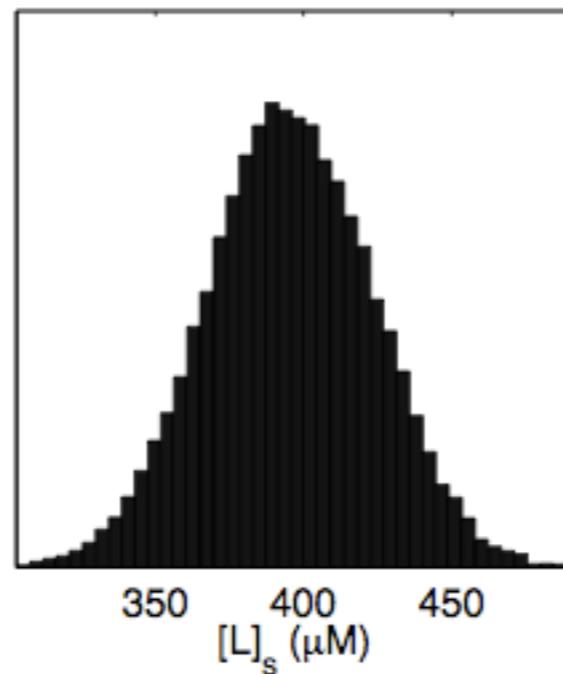
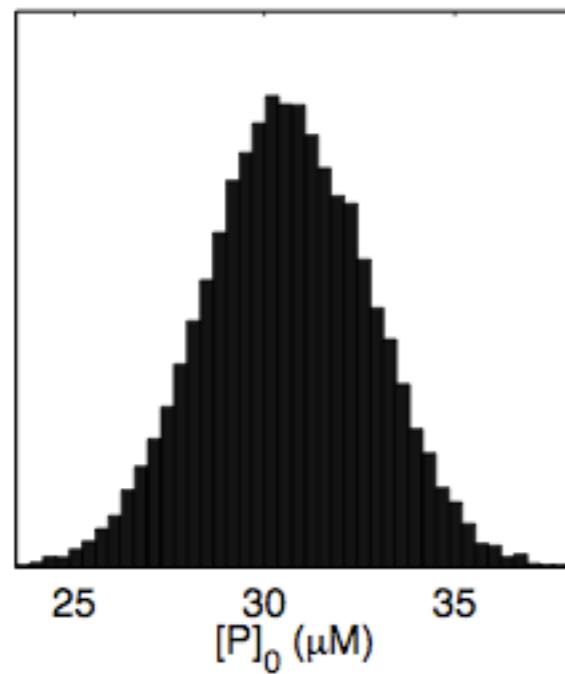
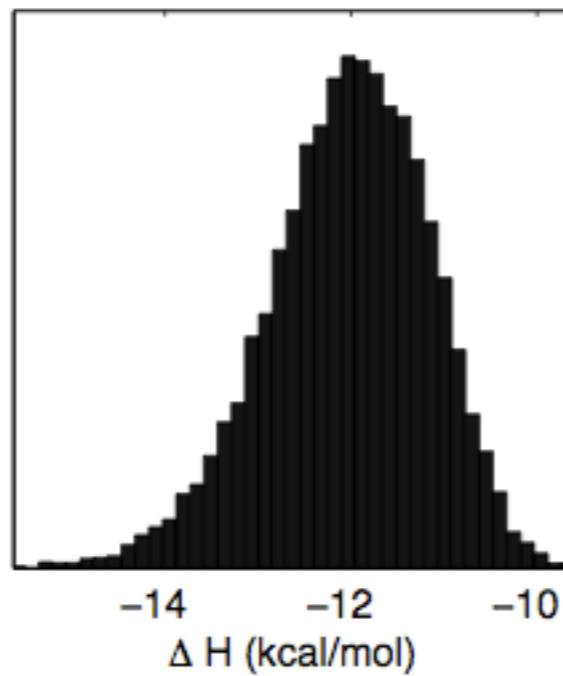
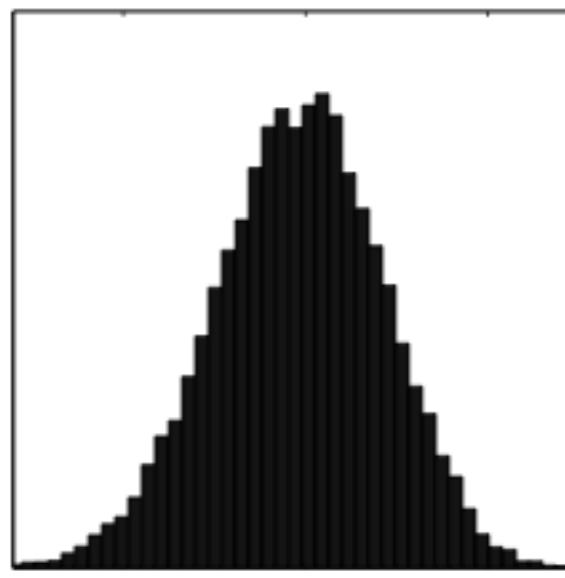
$p(\mathcal{D}|\theta)$ sampling distribution (model)

q_n^* nonlinear function of thermodynamic parameters

$p(\theta)$ prior

Analysis of ITC experiments: The Bayesian way

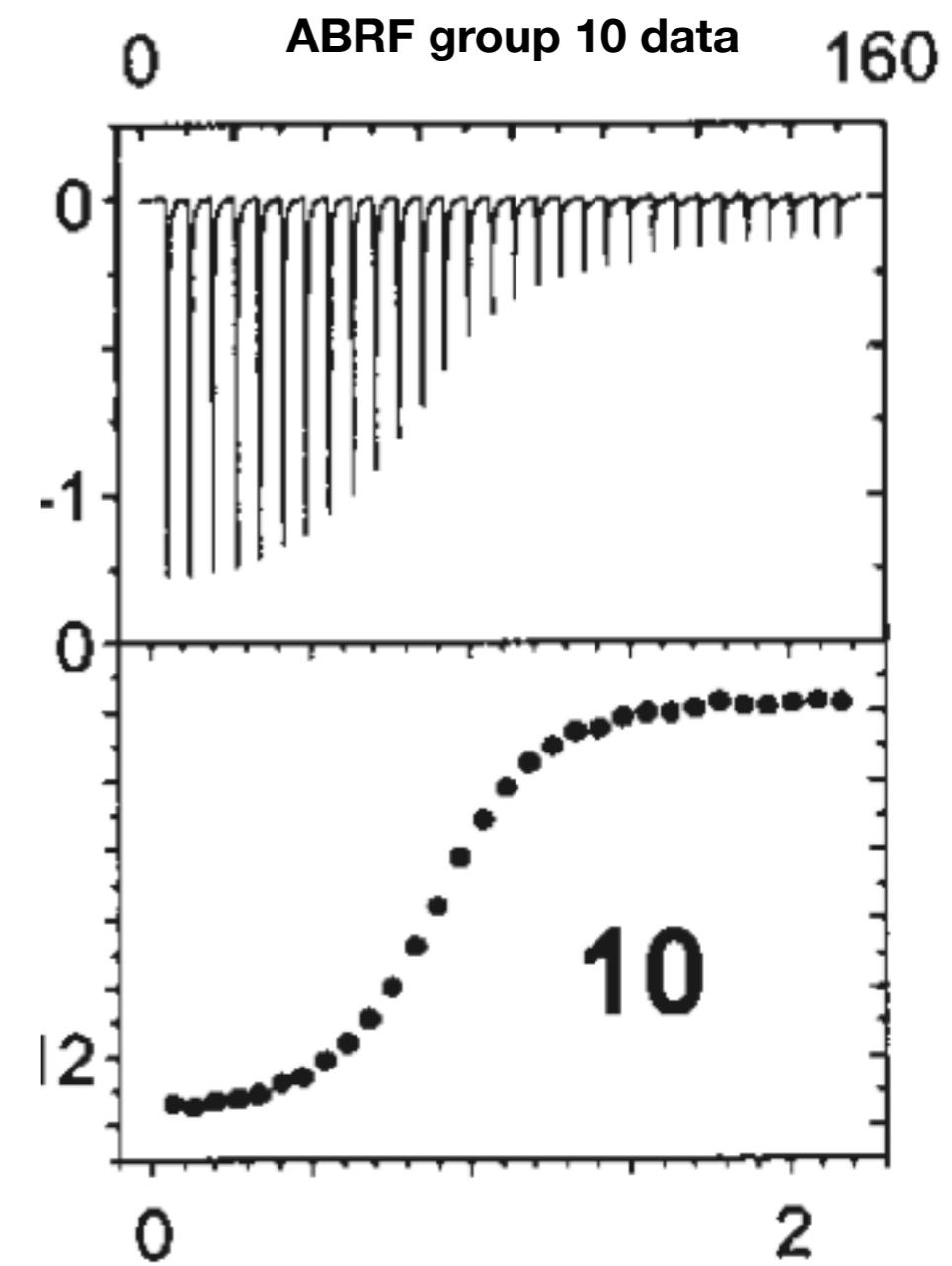
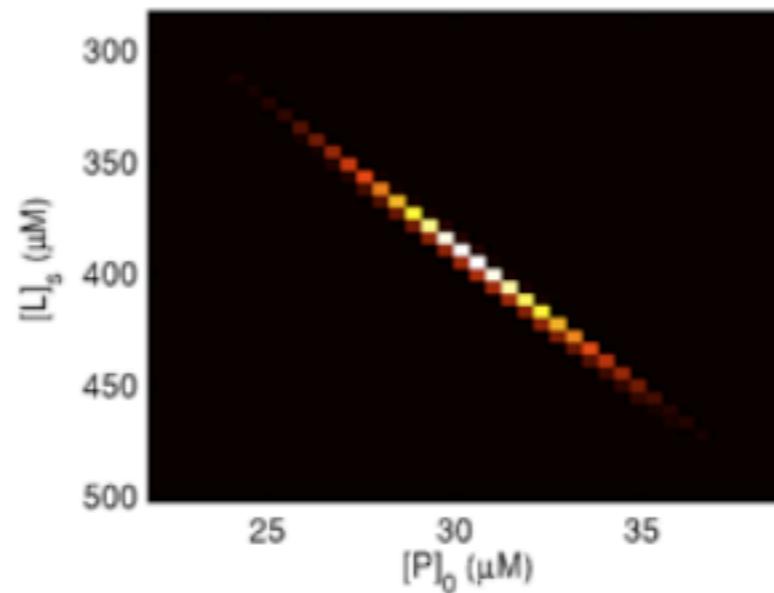
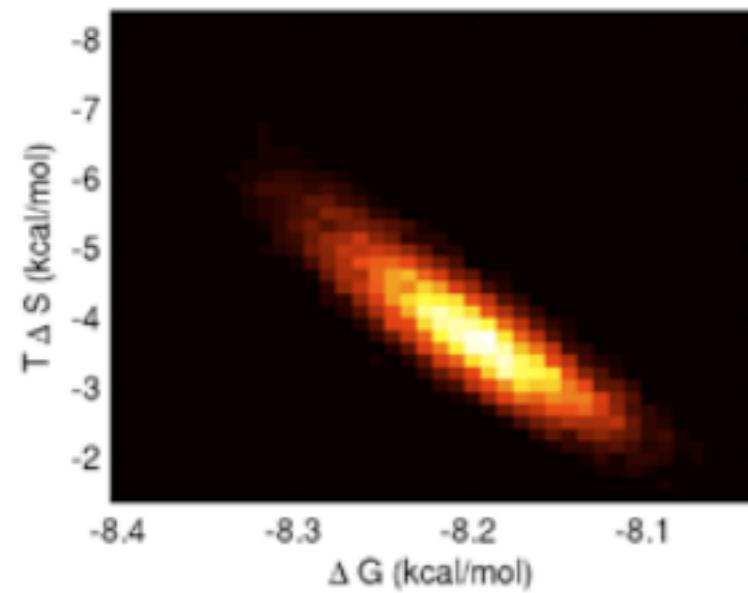
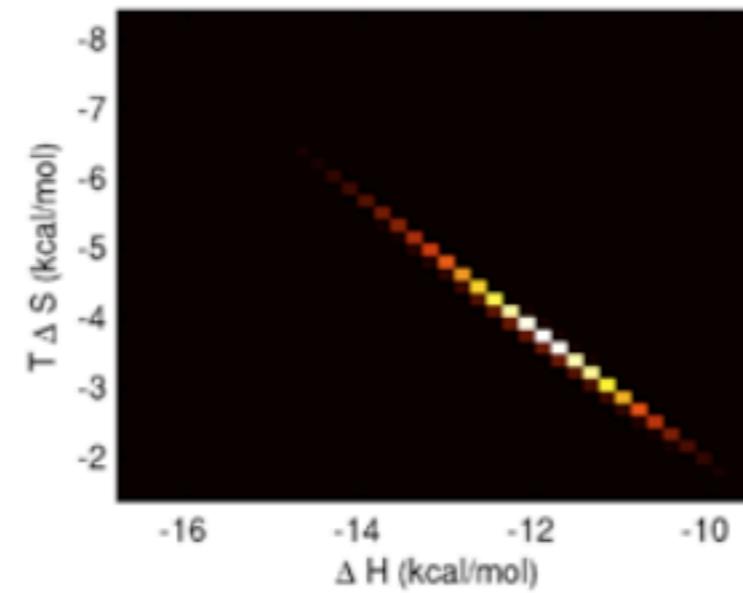
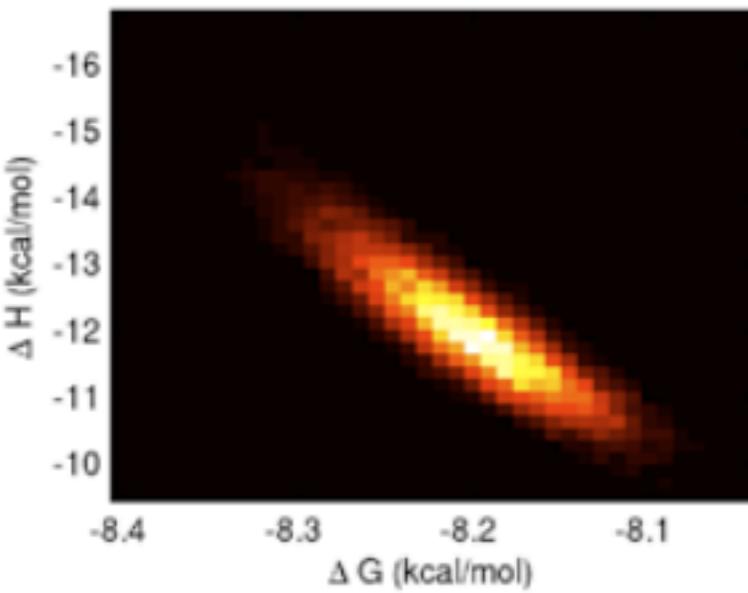
Marginal posterior probability distributions describe the uncertainty in each thermodynamic parameter inferred from the data.



Myszka et al. J. Biomol. Tech. 14:247, 2003.

Analysis of ITC experiments: The Bayesian way

Joint distribution functions describe the correlated uncertainty between any set of parameters.



Correlation in error in ΔH and $T\Delta S$ explains much of apparent entropy-enthalpy compensation behavior.

Myszka et al. J. Biomol. Tech. 14:247, 2003.

Bayesian approach has numerous advantages

Provides true posterior joint distribution of all thermodynamic parameters

Asymmetric confidence intervals and non-normal marginal distributions

Easy to “plug in” new binding models.

Can eliminate baseline “blank” experiment through titration in excess

Make joint inferences from data from multiple experiments

Instrument parameters can be conditioned on calibration data (e.g. NaCl titrations)

Expected information content of new experiments can be estimated for protocol design

Experimental design:

“Will experiment X give me enough information to make it worthwhile?”

“What is the best experimental design to reduce the uncertainty in Z?”

“Do I have to run a baseline for sample X?”

Code will be made available at <http://simtk.org/home/bayesian-itc>

Other explanations for apparent compensation? The “free energy window” effect: Instrumental limits?

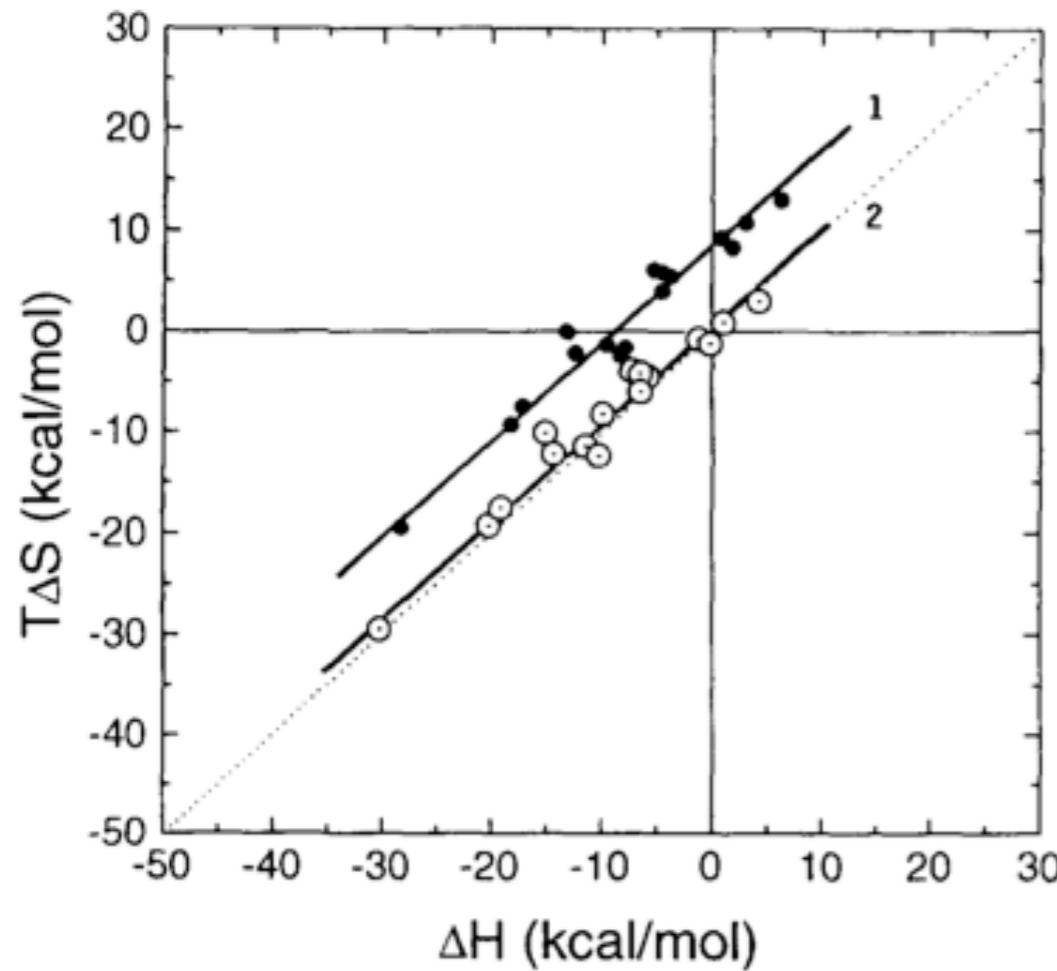
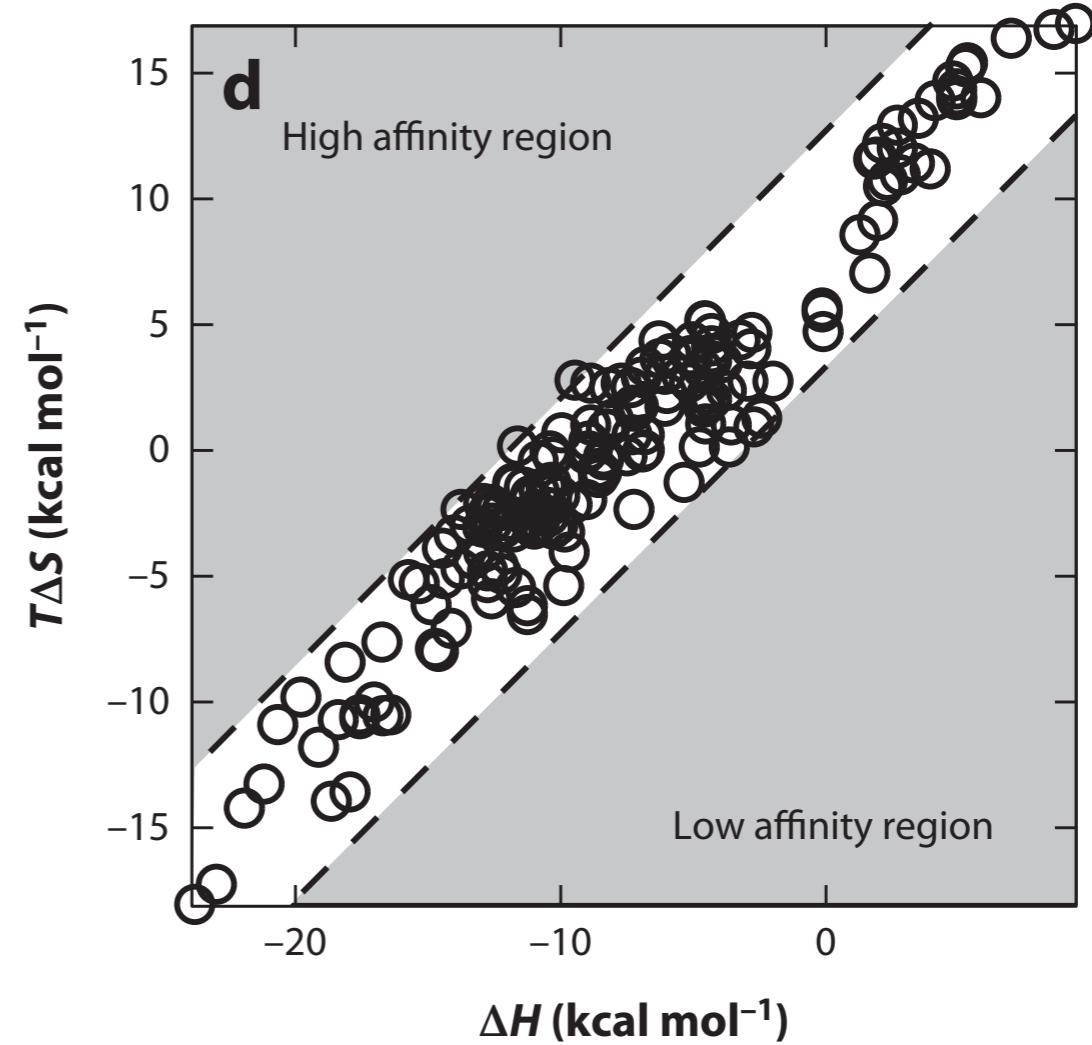


FIG. 3. The relationship between ΔH and $T\Delta S$ for the binding of Ca^{2+} to the proteins (●). The relationship between $\Delta H'$ and $T\Delta S'$ after subtracting the contribution of Ca^{2+} binding by itself (○). The data were taken from Table II. Lines 1 and 2 were obtained by a least squares fit of each data set. The dotted line represents $\Delta H - T\Delta S = 0$.



$$\Delta G = \Delta H - T\Delta S$$

All ITC measurements fall in a narrow range due to instrumental limitations

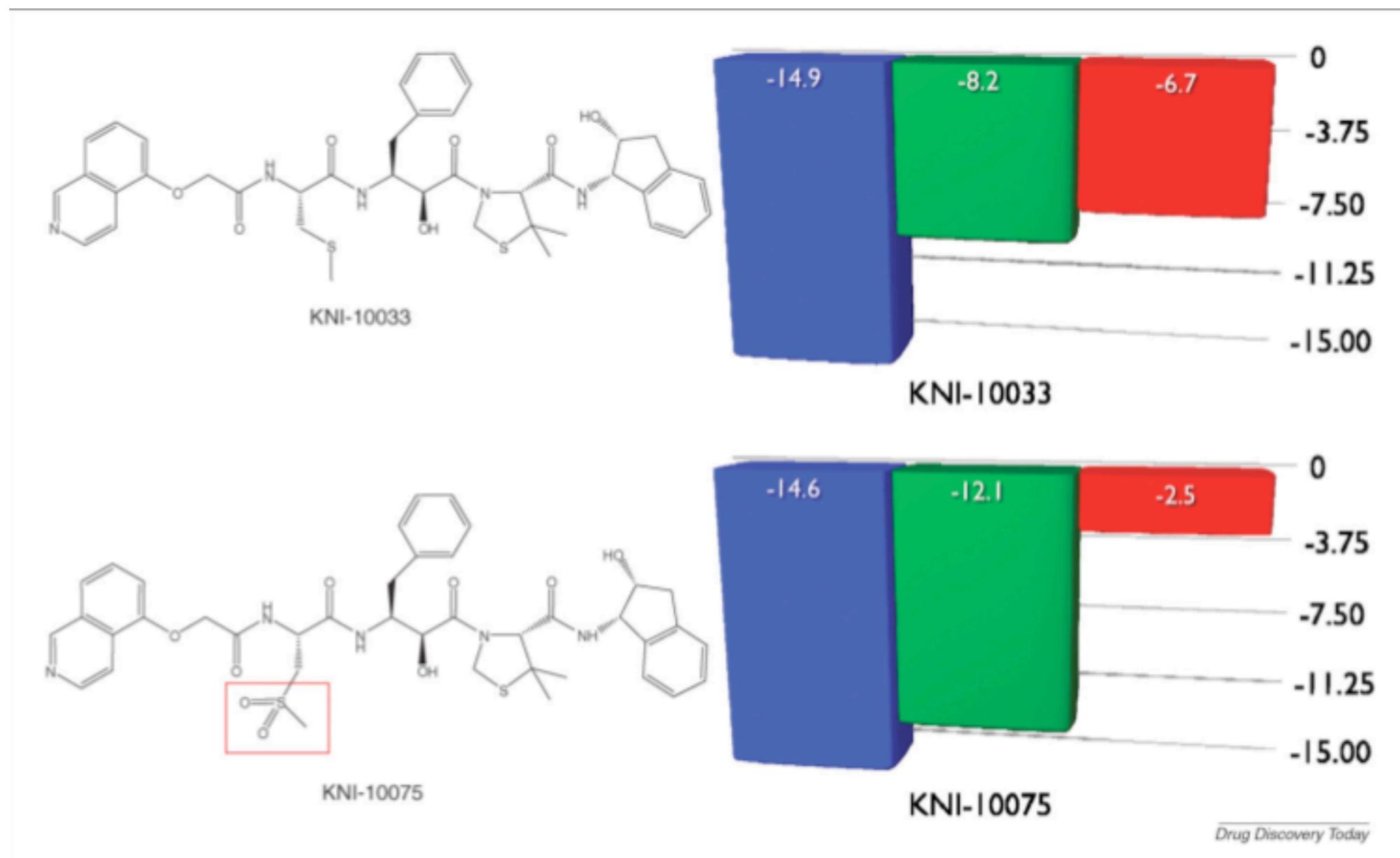
Exner O. 1997. How to get wrong results from good experimental data: a survey of incorrect applications of regression. *J. Phys. Org. Chem.* 10:797–813

Chodera and Mobley. *Annu Rev Biophys* 42:121, 2013.

Should we worry about compensation in drug design?

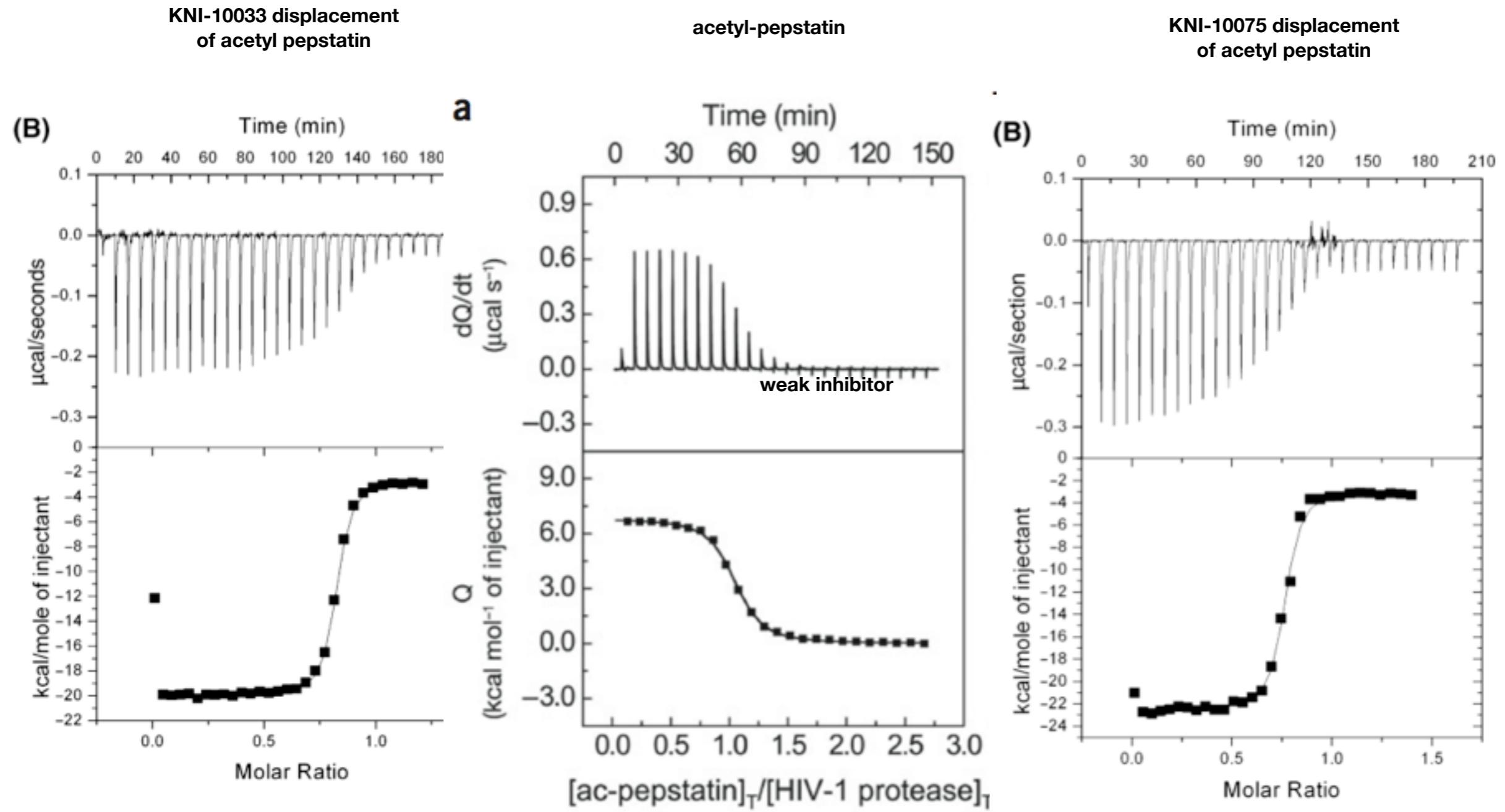
“The **thermodynamic signature**...provides a unique experimental way of characterizing the binding mode of a drug molecule.”

Freire E. *Drug Discovery Today* 13:869, 2009.



Does this claim hold up to what we've learned about correlated errors in ITC?

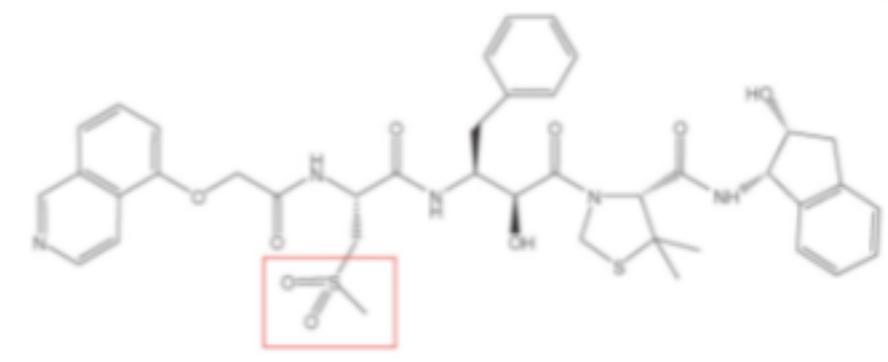
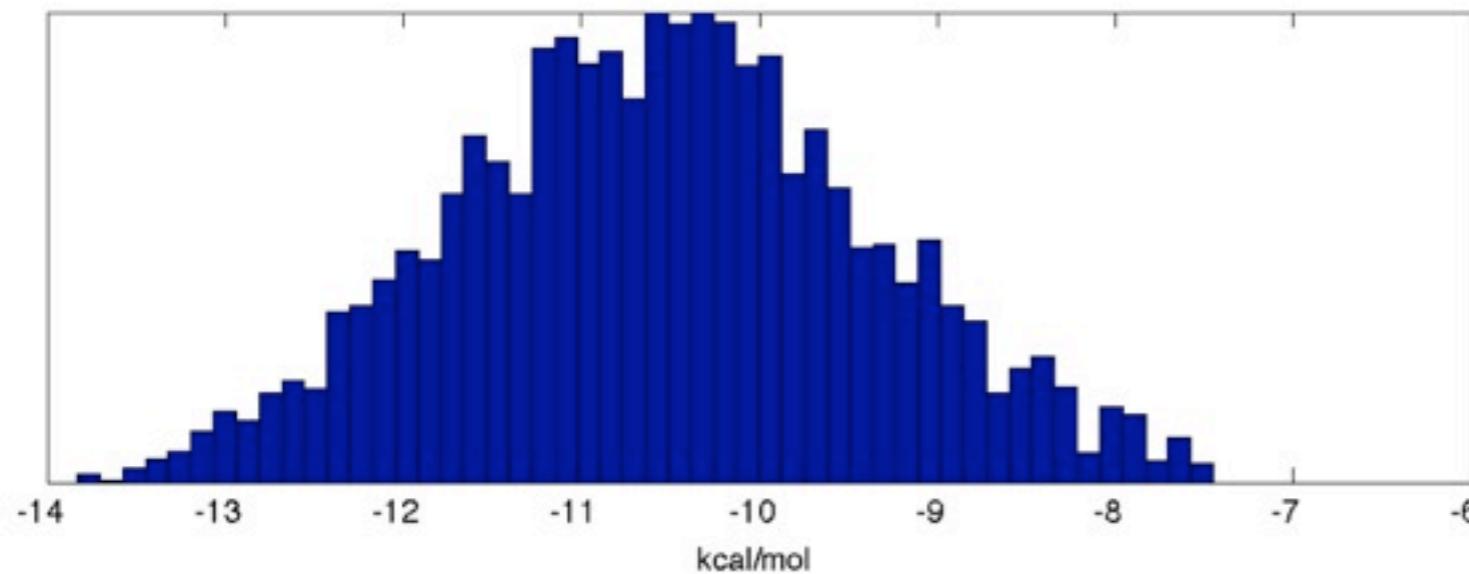
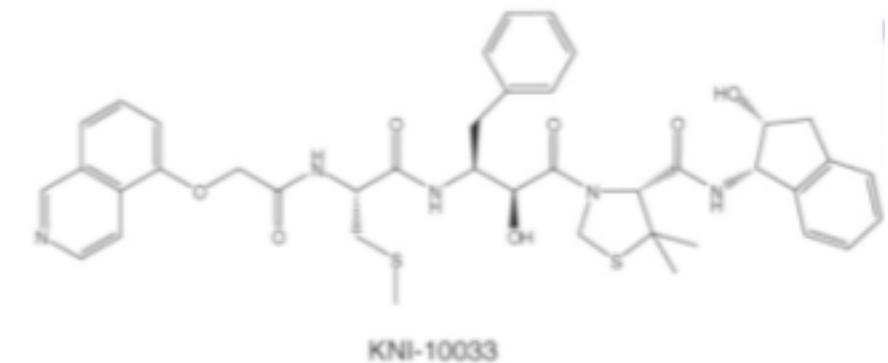
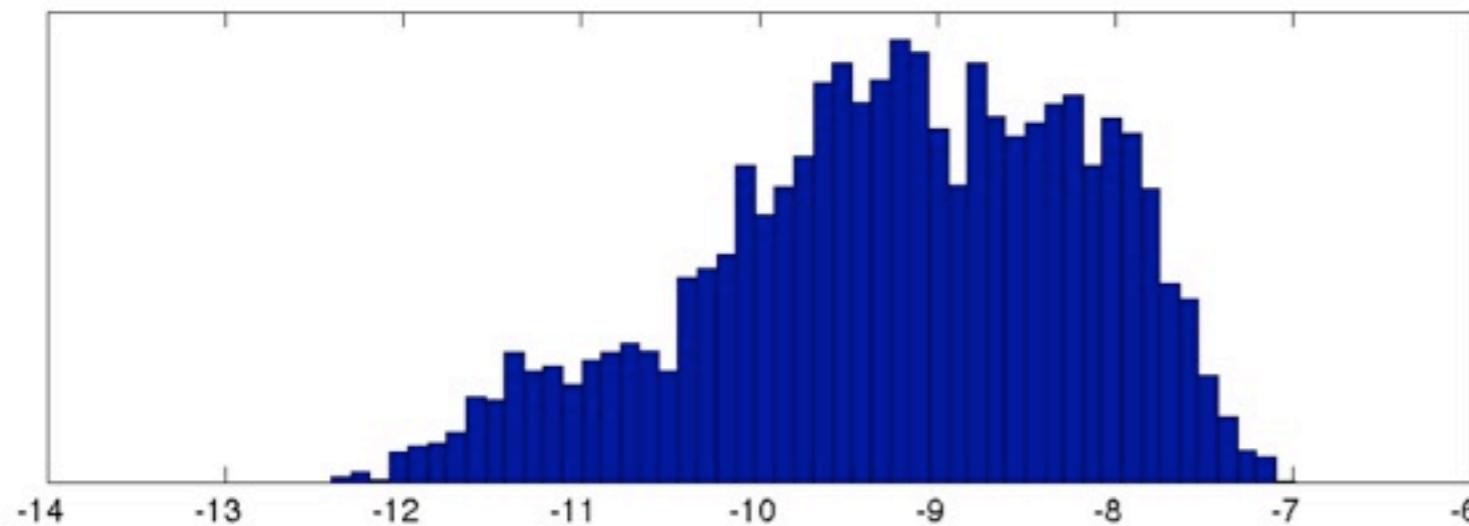
A more honest assessment of errors



Velazquez-Campoy A, Kiso Y, and Friere E. *Arch. Biochem. Biophys.* 390:169, 2001.
Velazquez-Campoy A and Friere E. *Nature Protocols* 1:186, 2006.

A more honest assessment of errors

ΔH of binding



Enthalpy change is clearly not the 4 kcal/mol difference the fingerprints would imply

Rules of thumb: Error magnitudes of last resort

Affinity/activity data:

Public Ki data: 0.54 pKi units

Kramer et al. J Med Chem 2012.

Public IC₅₀ data: 0.68 pIC₅₀ units

Kramer et al. PLoS One 8:e61007, 2013.

ITC data: 20% error in K_d and ΔH (unless evidence of propagating ligand concentration error)

Chodera and Mobley. Annu Rev Biophys 42:121, 2013; Myszka et al. J. Biomol. Tech. 14:247, 2003.

SPR data: 25% error in K_d, 20% in k_{on} and k_{off}

Myszka et al. J. Biomol. Tech. 14:247, 2003.

Concentrations:

prepared by a lab tech using mass/volume: ~10%

Myszka et al. J. Biomol. Tech. 14:247, 2003.

measured by NMR: 3%

Agilent technical note: Easy, Precise and Accurate Quantitative NMR [combining accuracy and precision]

An Editor's plea

Editorial policy, Author and Reviewer guidelines

Could we as a community gather what we've discussed this week and formulate it into clear, manageable, realistic guidelines that could be universally distributed, endorsed, and enforced?

Would this require some aspect of education or links to the same?

Corresponding Author Name: _____

Manuscript Number: _____

Reporting Checklist For Life Sciences Articles

This checklist is used to ensure good reporting standards and to improve the reproducibility of published results. For more information, please read [Reporting Life Sciences Research](#).

► Figure legends

Each figure legend should contain, for each panel where they are relevant:

- the **exact sample size (*n*)** for each experimental group/condition, given as a number, not a range;
- a **description of the sample collection** allowing the reader to understand whether the samples represent **technical or biological replicates** (including how many animals, litters, cultures, etc.);
- a **statement of how many times the experiment shown was replicated in the laboratory**;
- **definitions of statistical methods and measures**:
 - very common tests, such as *t*-test, simple χ^2 tests, Wilcoxon and Mann-Whitney tests, can be unambiguously identified by name only, but more complex techniques should be described in the methods section;
 - are tests one-sided or two-sided?
 - are there adjustments for multiple comparisons?
 - **statistical test results**, e.g., *P* values;
 - definition of ‘center values’ as **median or average**;
 - definition of **error bars as s.d. or s.e.m.**

Any descriptions too long for the figure legend should be included in the methods section.

Please ensure that the answers to the following questions are reported in the manuscript itself. We encourage you to include a specific subsection in the methods section for statistics, reagents and animal models. Below, provide the page number(s) or figure legend(s) where the information can be located.

► Statistics and general methods

Reported on page(s) or figure legend(s):

1. How was the sample size chosen to ensure adequate power to detect a pre-specified effect size?

2. Describe inclusion/exclusion criteria if samples or animals were excluded from the analysis. Were the criteria pre-established?

3. If a method of randomization was used to determine how samples/animals were allocated to experimental groups and processed, describe it.

For animal studies, include a statement about randomization even if no randomization was used.

4. If the investigator was blinded to the group allocation during the experiment and/or when assessing the outcome, state the extent of blinding.

For animal studies, include a statement about blinding even if no blinding was done.

5. For every figure, are statistical tests justified as appropriate?

Do the data meet the assumptions of the tests (e.g., normal distribution)?

Is there an estimate of variation within each group of data? Is the variance similar between the groups that are being statistically compared?

(Continues on following page)



► Reagents

Reported on page(s) or figure legend(s):

6. To show that antibodies were profiled for use in the system under study (assay and species), provide a citation, catalog number and/or clone number, supplementary information or reference to an antibody validation profile (e.g., [Antibodypedia](#), [1DegreeBio](#)).

7. Identify the source of cell lines and report if they were recently authenticated (e.g., by STR profiling) and tested for mycoplasma contamination.

► Animal models

Reported on page(s) or figure legend(s):

8. Report species, strain, sex and age of animals.

9. For experiments involving live vertebrates, include a statement of compliance with ethical regulations and identify the committee(s) approving the experiments.

10. We recommend consulting the ARRIVE guidelines (*PLoS Biol.* 8(6), e1000412, 2010) to ensure that other relevant aspects of animal studies are adequately reported.

► Human subjects

Reported on page(s) or figure legend(s):

11. Identify the committee(s) approving the study protocol.

12. Include a statement confirming that informed consent was obtained from all subjects.

13. For publication of patient photos, include a statement confirming that consent to publish was obtained.

14. Report the clinical trial registration number (at [ClinicalTrials.gov](#) or equivalent).

15. For phase II and III randomized controlled trials, please refer to the [CONSORT statement](#) and submit the CONSORT checklist with your submission.

16. For tumor marker prognostic studies, we recommend that you follow the [REMARK reporting guidelines](#).

► Data deposition

17. Provide accession codes for deposited data.

Reported on page(s) or figure legend(s):

Data deposition in a public repository is mandatory for:

- a. Protein, DNA and RNA sequences
- b. Macromolecular structures
- c. Crystallographic data for small molecules
- d. Microarray data

Deposition is strongly recommended for many other datasets for which structured public repositories exist; more details on our data policy are available [here](#). We encourage the provision of other source data in supplementary information or in unstructured repositories such as [Figshare](#) and [Dryad](#).

18. Is computer source code provided with the paper or deposited in a public repository? If so, indicate how it can be obtained.

References

- Christian Kramer and Richard Lewis, “QSARs, Data and Error in the Modern Age of Drug Discovery,” Current Topics in Medicinal Chemistry, 2012, 12, 1896-1902
- M. Hewitt, et al 2009, J. Chem. Inf. Model, 49,2572–2587

Acknowledgments

- All those experimentalists who are slaving in the trenches against the forces of nature and time to provide knowledge and insight to the ungrateful hoards.
- Especially those who have spent hours explaining to me what they do.
- Ant

More Acknowledgments

Ant (“I’m going to make you an offer you can’t refuse”) Nicholls

ITC and entropy-enthalpy compensation

**Kim Branson, Sarah Boyce, Paul Novick, Vijay Pande, David Minh,
David Mobley**

End of presentation

Terry Richard Stouch, PhD

President, Science for Solutions, LLC

Consulting in Drug Design; Pharmaceutical Research, Technologies, Process;
Molecular Simulation; Computational Sciences; Structural Biology

Senior Editor-in-Chief, Journal of Computer-Aided Molecular Design, Springer Publishing

AAAS Fellow

IUPAC Fellow

tstouch@gmail.com

1-609-275-7234