# Optimisation of Certainty-Based Assessment Scores

## A.R. Gardner-Medwin, Physiology (NPP), UCL, London

---

### Certainty-Based Marking (CBM)

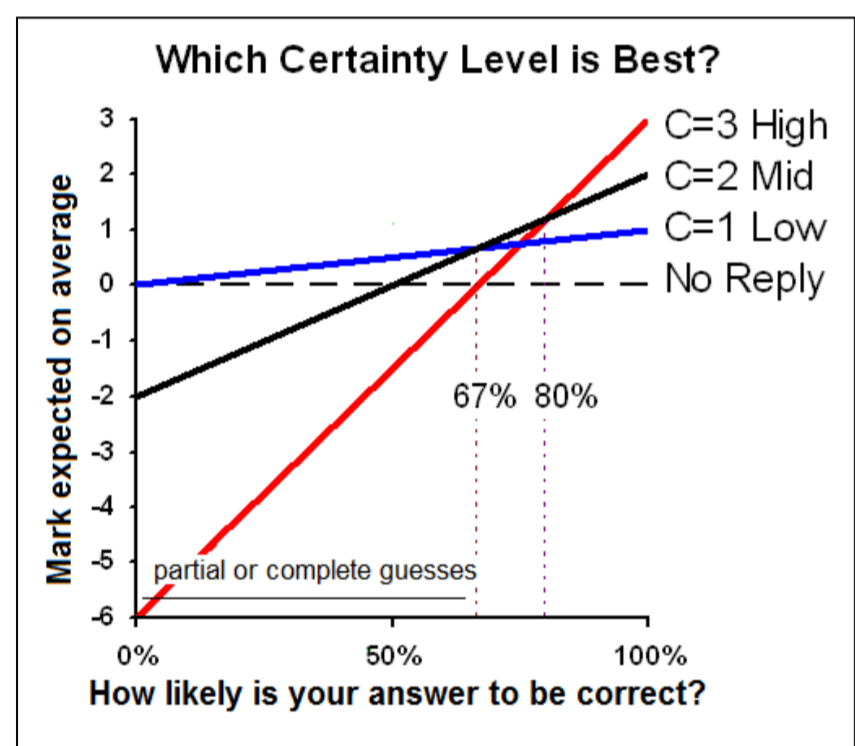**Aim :**  To optimise the presentation of
- Self-test Qs so as to challenge students and enhance study
- Exam Qs so as to increase the realism and predictive value of assessment data

### Background :
- Medicine and Physiology often require integration of knowledge from different perspectives to be sure of an answer.
- Thinking about the reliability of knowledge and inference is a key academic skill, with particularly dire consequences in Medicine when it fails.
- Valid measures of knowledge or ignorance must take account of uncertainty.
- Explicit certainty judgment has been shown in many psychological experiments to enhance learning and retention.

**CBM** is a *proper* mark scheme in the sense that a student is always motivated and rewarded for distinguishing and identifying honestly those answers that are uncertain and well justified. It is well founded in information theory (see THEORY, to the right).

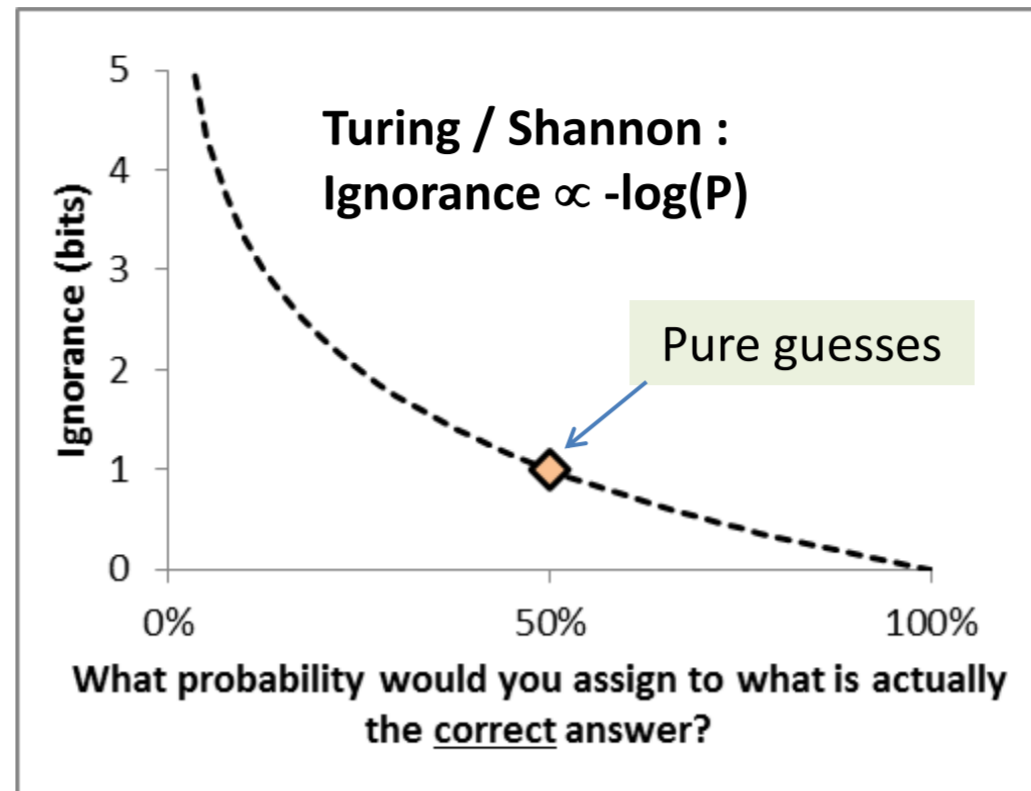| Degree of Certainty : | C=1 (low) | C=2 (mid) | C=3 (high) | No Reply |
|---|---|---|---|---|
| Mark if correct : | 1 | 2 | 3 | 0 |
| Penalty if wrong: | 0 | -2 | -6 | 0 |
| Probability Correct: | <67% | 67-80% | >80% | – |



**Which Certainty Level is Best?**

**Student perspective:**
- Always motivated to be honest
- Rewarded for identifying weaknesses
- Rewarded for sound justifications
- Encouraged to reflect & link info
- Misconceptions highlighted
- Simple and transparent scheme
- Perceive it as realistic & fair

**Staff perspective:**
- Doesn't require new or different Qs
- Enhanced feedback about content
- Enhanced reliability & validity in exams
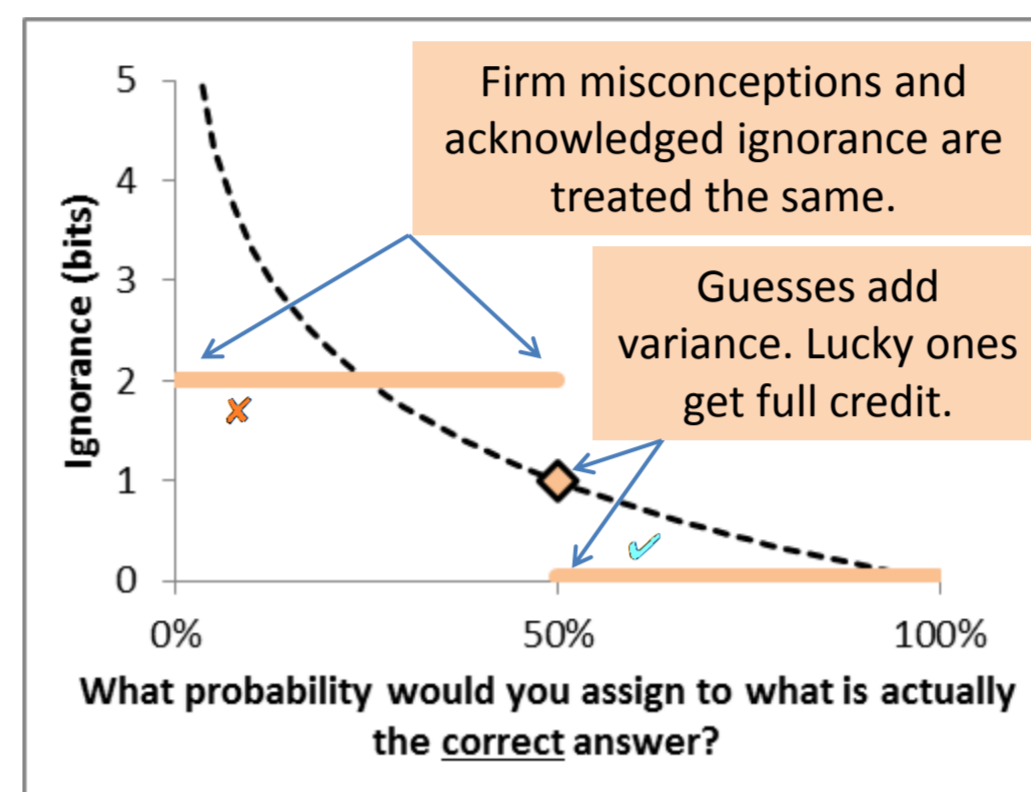- Better student learning experience

---

## AN  EXAMPLE

**Insulin injection raises blood glucose concentration.  True/False ?**

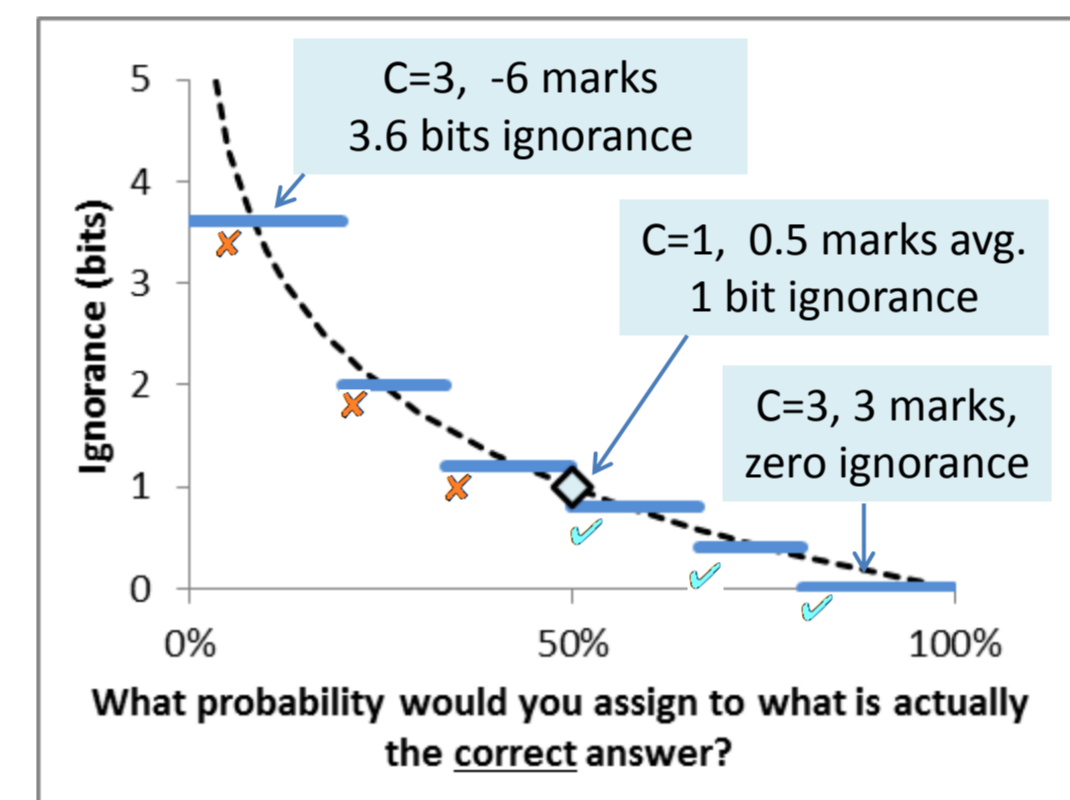### IGNORANCE  THEORY

*Ignorance has an unambiguous definition :*



Turing / Shannon : Ignorance ∝ -log(P)

Pure guesses

What probability would you assign to what is actually the **correct** answer?

*Ignorance inferred from binary (right / wrong) marks has serious problems:*



Firm misconceptions and acknowledged ignorance are treated the same.

Guesses add variance. Lucky ones get full credit.

What probability would you assign to what is actually the **correct** answer?

*(NB Negative marking generally make things worse: entering "Don't Know" instead of an uncertain answer will reduce variance, but will on average lose the student marks, unless the –ve marking is unusually severe)*

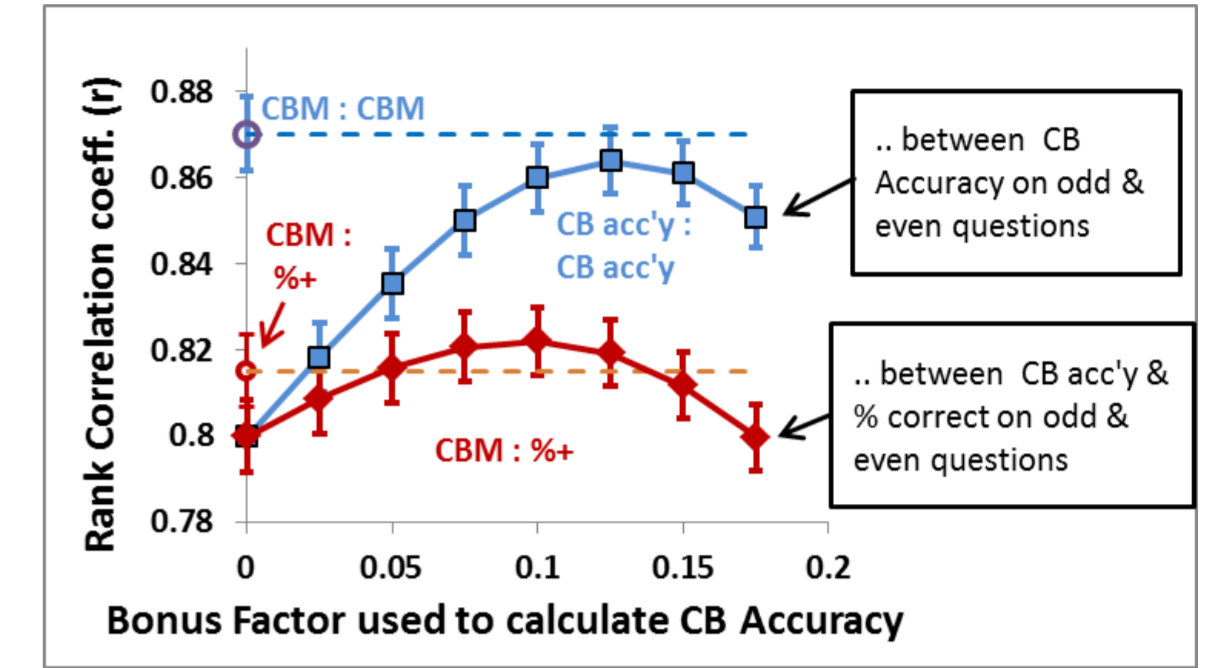*CBM gives a good measure of ignorance (= 0.4 x marks lost, relative to correct at C=3 ) :*



C=3, -6 marks 3.6 bits ignorance

C=1, 0.5 marks avg. 1 bit ignorance

C=3, 3 marks, zero ignorance

What probability would you assign to what is actually the **correct** answer?

*NB with Single-Best-Answer (SBA) Qs the situation is more complex, but the illustrated problems with conventional marking are more extreme.*
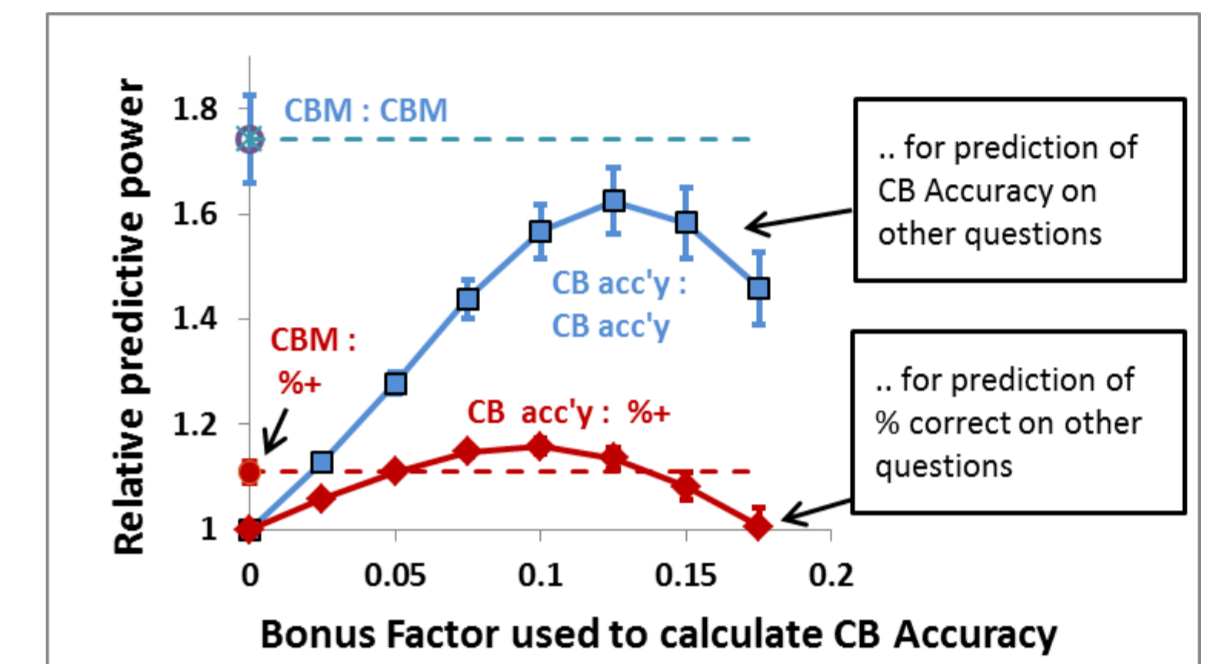
---

## Psychometric optimisation of CB Accuracy

A good measure of the quality of an assessment is how well the score or ranking based on one half of the test (e.g. odd numbered Qs) correlates with that based on the other half. This correlation is substantially enhanced with CBM. With "CB Accuracy" the scale of the "Bonus" added to Accuracy is a variable (<0.125 to ensure CBA<=100%). Data from 17 UCL Yr 1&2 medical exams are used here to assess the validity and reliabilty of CB Accuracy and to optimise the bonus factor (0.1 in the graphs presented in the left panel).

Mean correlation coeff ± s.e.m. for rankings based on odd & even no. Qs , using CBM totals or CB accuracy , paired with either the same score type  or with % correct.



.. between CB Accuracy on odd & even questions

.. between CB acc'y & % correct on odd & even questions

Enhancement of predictive power of answers (increase of r/(1-r) compared with right/wrong marking) based on the above correlations. A factor of 1.5 means that the benefit is approximately equivalent to a 50% increase in Q numbers in a conventional test.



.. for prediction of CB Accuracy on other questions

.. for prediction of % correct on other questions

A bonus factor of 0.1 is clearly a good choice, giving nearly as good psychometric reliability as the CBM mark itself, and better prediction of accuracy and rank based on accuracy on the complementary Q set.

---

## CBM Implementation: www.ucl.ac.uk/LAPT

All you need to implement CBM for self-tests in your institution (following a model developed for Imperial & Kings, London), are:
- A server site where students & staff authenticate with a local userid
- Links to that site, specifying each self-test, in your VLE

You can use open exercise files, or private ones sited either on your server or at UCL. Editing is simple. Contact me at ucgbarg@ucl.ac.uk for more information. Wholly self-contained software packages are under development, but server loading for new users is almost negligible because computation nearly all takes place on the student's computer.

---

## SUMMARY

**CBM makes sense!**

**Doesn't require special Questions**

**Always motivates students to give careful honest judgment**

**SELF-TESTS**
- More sound and fair measure
- ↑ reflection & linking of Info
- ↑ realism about uncertainty
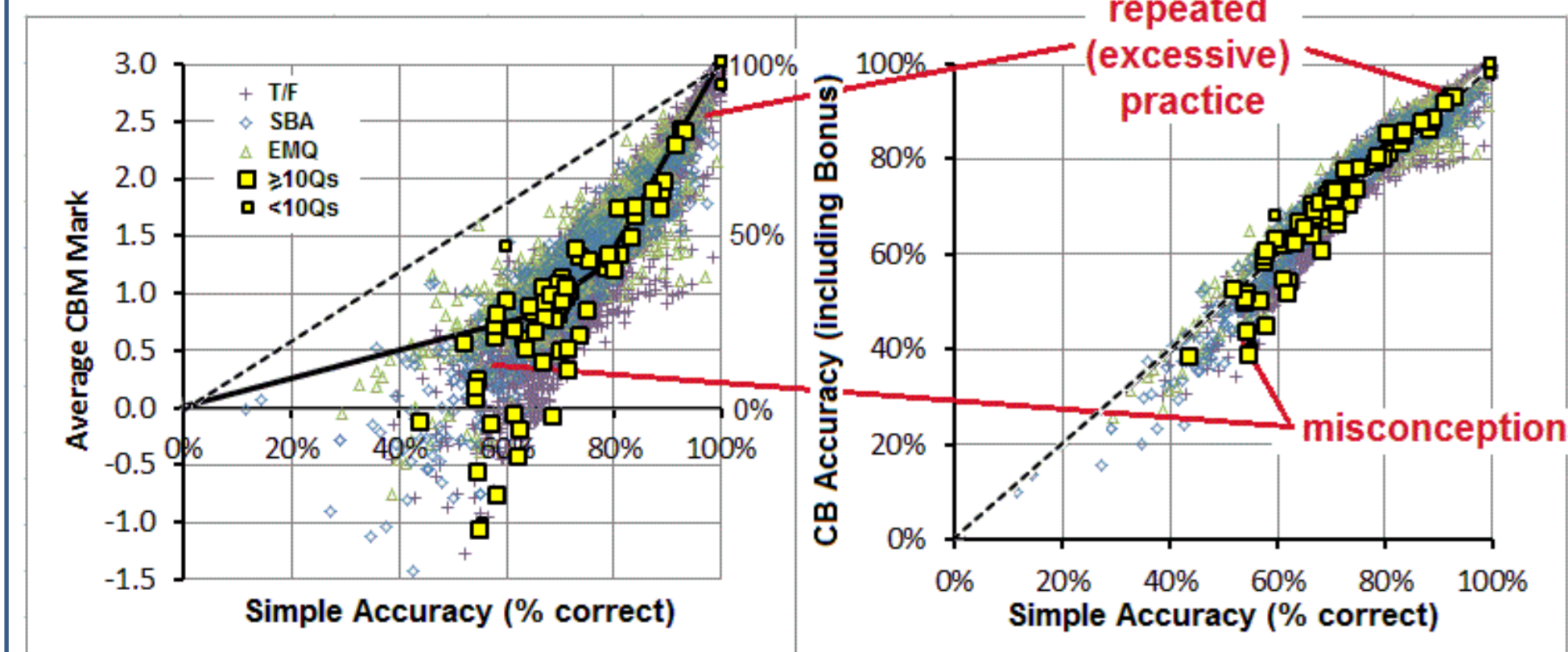- Highlights misconceptions
- Students like it!

**EXAMS**
- ↑ psychometric reliability
- ↑ psychometric validity
- ↓ question numbers
- No loss of conventional exam info

---

### Some Distinguished Tweets !

"When you know a thing, to hold that you know it, when you do not know a thing, to allow that you do not know it – this is knowledge."
*Confucius*

"… there are known knowns; ... there are known unknowns; ... But there are also unknown unknowns"
*Donald Rumsfeld*

"It's not ignorance does so much damage; - it's knowin' so derned much that ain't so."
*attr.: Josh Billings*

"A lucky guess is not  knowledge. A firm misconception is worse than acknowledged ignorance. So why do we mark students as if these things weren't true?"
*TG-M*
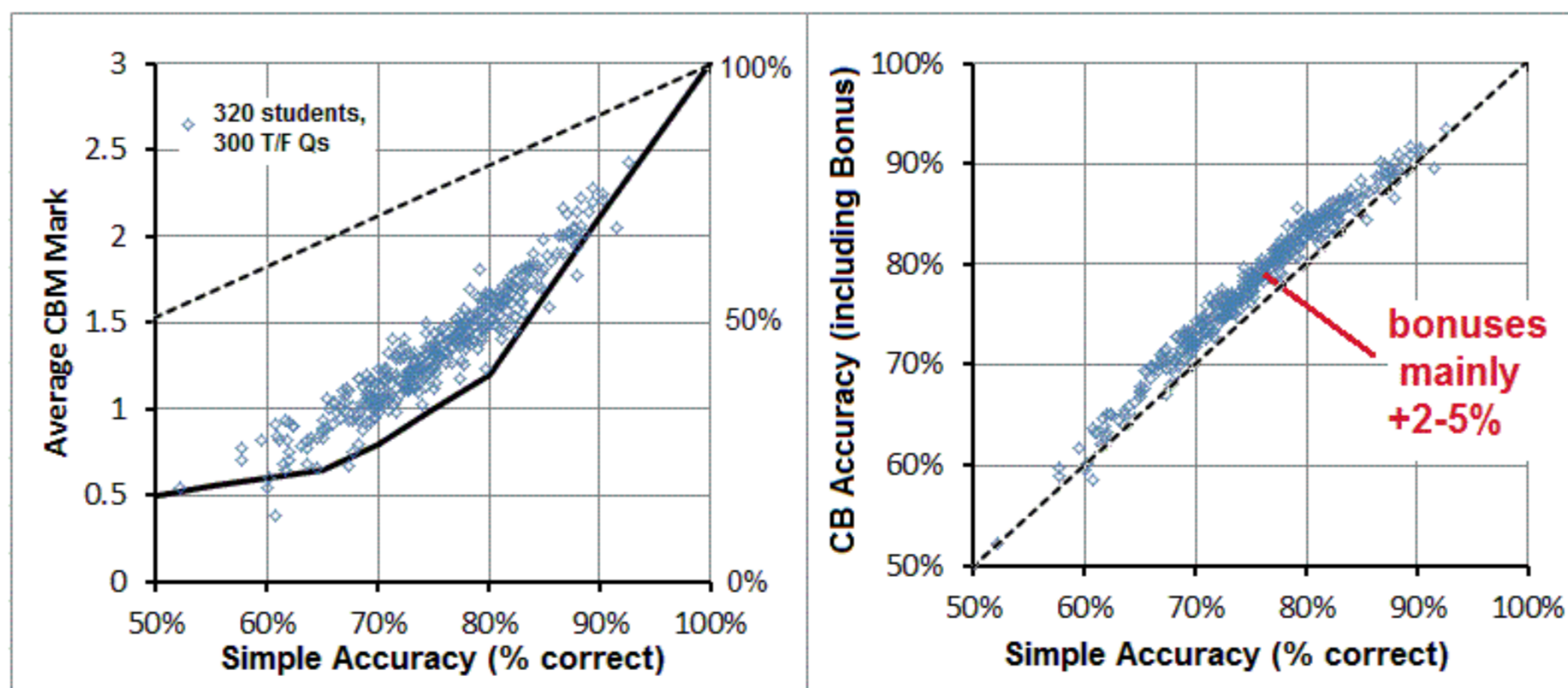
---

## HOW BEST TO PRESENT CBM SCORES

CBM motivates a student to reflect and identify uncertain vs. reliable answers. This is how you maximise your score. But no student can realistically expect to get an average CB mark as great (expressed as a % of maximum) as their accuracy, or % correct on a test.  This would only be attainable if every correct answer was given at C=3 and every error at C=1, which in most tests is unrealistic. The graphs below show how (both in self-tests and exams) 80% accuracy is typically associated with an average CBM mark = 1.5, only 50% of maximum.

**THE PROBLEM!** This simple comparison, even for students above average at judging uncertainty, can be demoralising and counterproductive. There is nothing wrong with CBM scores; but they are fundamentally different from (and psychometrically superior to) accuracy measures. The problem of presentation is tackled here by generating a **"CB Accuracy"** by adding a **BONUS** to the simple accuracy as a measure of how well the student categorises responses as uncertain or reliable. The bonus is positive or negative, proportional to the amount the average CBM mark is above (or below) the average that would be obtained (shown by heavy black lines below) if the student had used the same optimal C level for all his/her answers. Negative bonuses are common in self-tests when students often have misconceptions (confident errors), but as is evident in exam data, students can aspire to gaining positive bonuses of 2-5%.
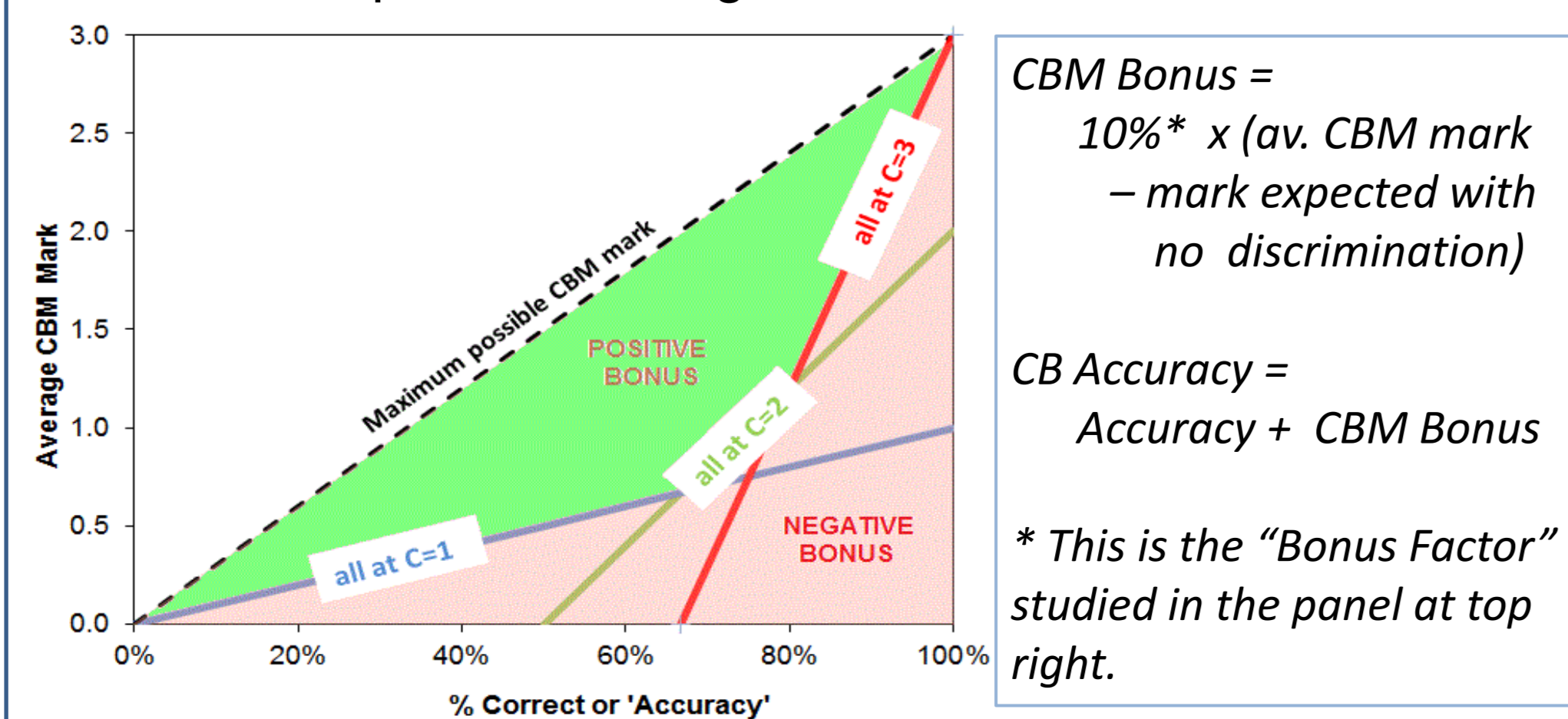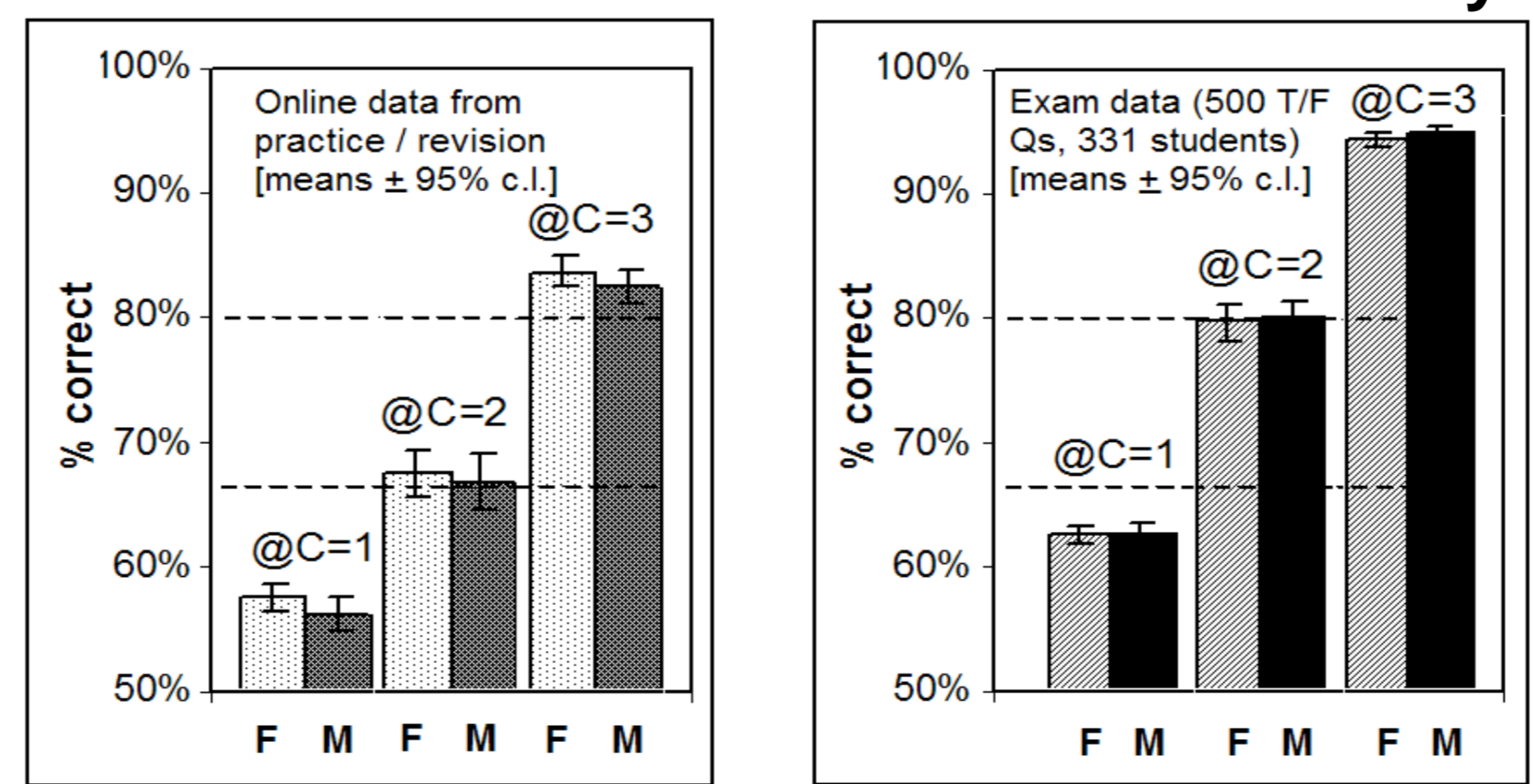
**Online self-test practice data**



repeated (excessive) practice

misconceptions

**Exam Data**



bonuses mainly +2-5%

Yellow squares show students' scores on self-test physiology T/F Qs, as displayed in feedback to students and staff using LAPT (www.ucl.ac.uk/LAPT). They are superimposed for comparison on a background of historic data from ca. 9000 sessions with various Q types. Exam data is from a 1st yr medical exam at UCL after students had had substantial experience using CBM in LAPT.
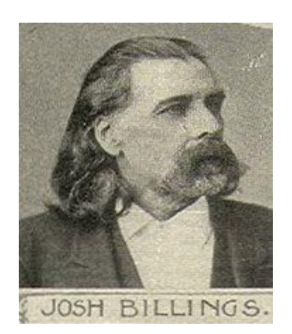
### How well do students discriminate reliability?



Online data from practice / revision [means ± 95% c.i.]

Exam data (500 T/F Qs, 331 students) [means ± 95% c.i.]



CBM Bonus = 10%* x (av. CBM mark – mark expected with no  discrimination)

CB Accuracy = Accuracy +  CBM Bonus

* This is the "Bonus Factor" studied in the panel at top right.

---