

# STATISTICAL PRACTICE AT THE BELLE EXPERIMENT, AND SOME QUESTIONS

*Bruce Yabsley*

Virginia Polytechnic Institute and State University, Blacksburg VA

## Abstract

The Belle collaboration operates a general-purpose detector at the KEKB asymmetric energy  $e^+e^-$  collider, performing a wide range of measurements in beauty, charm, tau and 2-photon physics. In this paper, the treatment of statistical problems in past and present Belle measurements is reviewed. Some open questions, such as the preferred method for quoting rare decay results, and the statistical treatment of the new  $B^0/\bar{B}^0 \rightarrow \pi^+\pi^-$  analysis, are discussed.

## 1 INTRODUCTION

My ambitions for this conference are to recover my luggage, which went missing four days ago, and to find answers to some questions about statistical practice at Belle.<sup>1</sup> I suspect I'm not the only one with an agenda in this area. From the point of view of the Belle spokesmen, it would be far better if I could articulate workable guidelines for our use of statistical methods . . . rather than finding fault with our present practice case-by-case. As I hope to show, it is much easier to criticise our statistical practice than it is to suggest alternatives, although I have made some tentative steps in that direction.

Analyses at Belle are not all of the same kind, and the “statistical environment” varies from one study to another. After reviewing the experiment itself (section 2) and some general issues (section 3), I hope to give you a feeling for the main types of analysis, and the statistical issues in each case (section 4). Of particular interest are the searches for rare B-meson decays (section 4.2), where there is a tradeoff between “purist” statistical concerns and important practical issues; and our new analysis of  $B^0/\bar{B}^0 \rightarrow \pi^+\pi^-$  decays (section 5). Interpretation of this result is surprisingly difficult, due to the unusual configuration of the parameter space, as well as some features of the analysis. The first Belle paper on this topic is two weeks from journal submission, so you may have the chance to influence our statistical practice in real time!

## 2 THE BELLE EXPERIMENT

The main “point” of Belle is to test the Kobayashi-Maskawa model of CP-violation [1], in which the phenomenon is entirely due to the irreducible complex phase of the quark mixing matrix. The simplicity of this model makes it highly predictive. In the absence of extra generations or isosinglet quarks, we have the familiar  $3 \times 3$  matrix

$$\begin{pmatrix} V_{ud} & V_{us} & V_{ub} \\ V_{cd} & V_{cs} & V_{cb} \\ V_{td} & V_{ts} & V_{tb} \end{pmatrix}$$

(the so-called CKM matrix), and unitarity gives a number of relations of the form  $V_{td}V_{tb}^* + V_{cd}V_{cb}^* + V_{ud}V_{ub}^* = 0$ . This particular “unitarity triangle” (see Fig. 1) is especially interesting: all its interior angles  $\phi_{1,2,3}$  are believed to be far from  $0^\circ$  and  $180^\circ$ , and may be measured using time-dependent asymmetries between  $B^0$  and  $\bar{B}^0$  decays to appropriate states. In other words: CP-violating asymmetries in these decays are expected to be large. The recent measurement of  $\sin 2\phi_1$  by Belle [2] confirms that this

---

<sup>1</sup>My luggage arrived on the Thursday morning, about 30 hours before I left Durham. Some of my questions were answered, at least in part. More of this anon.

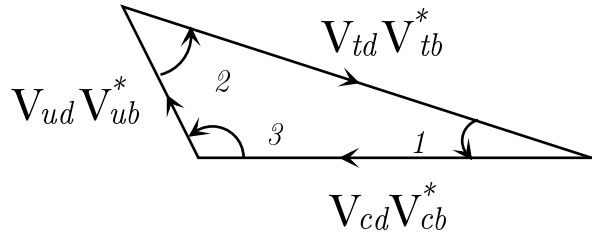


Fig. 1: The unitarity triangle relevant to B decays. In all unitarity triangles, the angles differ from  $0^\circ$  and  $180^\circ$  if the KM phase is nonvanishing. In *this* triangle, the angles are expected to be far from  $0^\circ$  and  $180^\circ$ , leading to readily measurable effects.

is true, at least regarding  $\phi_1$ : based on a study of  $B^0, \bar{B}^0 \rightarrow J/\psi K_S^0$  and related decays, we found  $\sin 2\phi_1 = 0.99 \pm 0.14 \pm 0.06$ . (This is the first observation of a CP-violating effect outside the neutral kaon system. In a spirit of friendliness let me also cite the similar, and essentially simultaneous, result from BaBar [3].) The  $B^0/\bar{B}^0 \rightarrow \pi^+\pi^-$  analysis discussed in section 5 is sensitive to the angle  $\phi_2$ .

The Belle detector [4] is a classic solenoidal tracking detector ( $|\vec{B}| = 1.5\text{ T}$ ) with some extra features. Most notable is its asymmetry, to optimize acceptance of events with a centre-of-mass boosted by  $\beta\gamma = 0.425$  in the lab frame. (The storage rings of KEKB collide  $8\text{ GeV } e^-$  on  $3.5\text{ GeV } e^+$ , *i.e.*  $\sqrt{s} = 10.58\text{ GeV}$ , yielding  $\Upsilon(4S) \rightarrow B\bar{B}$  decays or nonresonant  $q\bar{q}, \tau^+\tau^-$  *etc.*) The drift of the  $B^0/\bar{B}^0$  prior to decay is measured using a silicon vertex detector, allowing measurement of time-dependent asymmetries (see section 4.1). Particle ID over the full momentum range for B-daughters is obtained using aerogel Čerenkov counters as well as the more traditional  $dE/dx$  and time-of-flight measurements: this improves B flavour-tagging, and allows identification of rare B decays where we hope to find evidence of “direct” CP violation (section 4.2). The flux return is instrumented with RPCs to identify  $\mu^\pm$  (as matched tracks), improving the efficiency and purity of  $J/\psi \rightarrow \mu^+\mu^-$  reconstruction, and  $K_L^0$  (as neutral clusters), allowing measurement of  $\sin 2\phi_1$  using  $J/\psi K_L^0$  and related modes.

The collaboration itself is of order 250 physicists, from over 50 institutes in 12 countries. Perhaps half of the analysis effort is directed towards headline studies like  $J/\psi K^{(*)}$  and  $\pi^+\pi^-$ , with the rest thinly spread over rare B decays, CKM matrix elements (*via* measurements of semileptonic decays), charm studies,  $\tau$  and  $2\gamma$  physics. Most of this work proceeds through some mix of local “specialties” and individual initiative, with fairly light coordination from the centre. Which brings us to my next point.

### 3 ABOUT STATISTICAL PRACTICE AT BELLE

*In those days there was no king in Israel; everyone did what was right in his own eyes.*

Judges 21:25

The key to understanding statistical practice at Belle is that *there is no official policy*. Partly as a result of this, our practice is inconsistent. The following is my own take on this state of affairs.

- This is not nearly as bad as it sounds:
  - Most of our procedures *are* motivated, either by principle or tradition. In some places traditional methods have perhaps been mistaken for rigorous formulae, but this is a higher-order problem: one of education, not discipline.
  - We do describe our procedures in our papers, and this is more important than consistency. The description is sometimes incomplete, but I hope we are becoming more sensitive to this.
- There *are* things which should, and may change. I will mention some of them below.
- In avoiding anarchy, we do not want to become an authoritarian state.

The last of these points is worth stressing. Most of us spend our time going about our own business, and many of us have significant freedom to *define* what that business is. In part this is our way of working:

250 physicists may sound like a lot, but the number of potential physics topics is vast, and individuals “fanning out” opportunistically is a good way to cover them. But it is a *modus vivendi* as well. We live with each others’ egos, and satisfy our own, by spending most of our time on our own affairs.

Prescriptive policies have a way of cutting across this, and we tend to avoid them. Of course we set standards for papers (our *ad hoc* paper refereeing committees are almost the only committees worthy of the name) and areas of analysis are run by conveners, who act as facilitators and clearinghouses for information; occasionally, we also commission work. But to write a recipe-book for statistical practice, which all published analyses were expected to follow—this would be unprecedented. I don’t think people would be happy with it, and I doubt “the authorities” would exert themselves to enforce it . . . change, if it occurs, is more likely to proceed *via* the winning of hearts and minds. I suspect consciousness-raising in analysis groups, and among the individuals who sit on those refereeing committees, is the key.<sup>2</sup>

## 4 THE BELLE ANALYSES

The analyses themselves, by which I mean analyses we have written up and published (or have under review), can be divided into four broad categories:

1. the flagship time-distribution fit analyses (usually concerning CPV);
2. searches for rare hadronic  $B$  decays;
3. fits over a flattish background;
4. systematic-dominated analyses.

This is a grouping according to statistical problems, rather than common work. The second category corresponds roughly to one of our analysis groups, which does lead to a certain uniformity of method. Otherwise, these categories cut across Belle’s administrative divisions. I will consider them in turn, with the aid of an example in each case.

### 4.1 The flagship time-distribution fits

The unitarity angle measurements mentioned above, based on  $B^0/\bar{B}^0 \rightarrow J/\psi K^{(*)}$  *etc.* [2] and  $\pi^+\pi^-$  [5], are performed by fitting decay-time distributions. To illustrate the method I have chosen a simpler analysis, which (for our purpose) has the same features: a measurement of the  $D^0 - \bar{D}^0$  mixing parameter  $y_{CP}$ , through the lifetime difference of neutral D-mesons decaying to  $K^-K^+$  (a CP eigenstate) and  $K^-\pi^+$  [6].

After selecting  $D^0 \rightarrow K^-K^+$  and  $K^-\pi^+$  decays<sup>3</sup> (imposing particle ID and decay angle cuts, and a momentum cut to veto B-daughters) we fit the tracks to a common vertex, and extrapolate the resulting D-momentum  $\vec{p}_D$  to the interaction region: this gives us the flight length and thus the decay time. Thanks to our good kaon/pion separation the D-decay samples are fairly clean, so we avoid  $D^*(2010)^+ \rightarrow D^0\pi^+$  tagging, keeping the samples as large as possible. Of course some combinatorial  $K^-K^+/K^-\pi^+$  is accepted, and we estimate the probability for each candidate  $i$  to be a true D-decay using

$$f_{SIG}^i = \frac{N_{SIG}(M^i)}{(N_{SIG}(M^i) + N_{BKG}(M^i))}, \quad (1)$$

where  $M^i$  is the mass of the candidate, and signal and background fractions  $N_{SIG}(M^i)$  and  $N_{BKG}(M^i)$  are taken from a fit to the mass distribution of all candidates. The distributions are uncomplicated, and double-Gaussian fits over linear backgrounds are sufficient for the purpose. We accept events well into the tail—out to  $6\sigma$  in the mass—for reasons that should become clear.

<sup>2</sup>When I began preparing this study I thought a set of usable public tools was the key. While these have their place, and I haven’t abandoned the ambition of providing some, I’ve come to believe that issues of principle, and some genuinely open questions, must be addressed first. The remainder of this talk hopefully shows why.

<sup>3</sup>Here and in general, I imply the inclusion of charge-conjugate modes.

We then perform an unbinned maximum likelihood fit to the events, using the function

$$\begin{aligned} \mathcal{L}(\tau_{SIG}, S, S_{tail}, f_{tail}, \tau_{BG}, f_{\tau_{BG}}, S_{BG}, S_{tail}^{BG}, f_{tail}^{BG}) \\ = \prod_i \left[ f_{SIG}^i \int_0^\infty dt' \frac{1}{\tau_{SIG}} e^{-\frac{t'}{\tau_{SIG}}} R(t^i - t'; \sigma_t^i, S, S_{tail}, f_{tail}) \right. \\ \left. + (1 - f_{SIG}^i) \int_0^\infty dt' \left\{ f_{\tau_{BG}} \frac{1}{\tau_{BG}} e^{-\frac{t'}{\tau_{BG}}} + (1 - f_{\tau_{BG}}) \delta(t') \right\} R(t^i - t'; \sigma_t^i, S_{BG}, S_{tail}^{BG}, f_{tail}^{BG}) \right] \end{aligned} \quad (2)$$

which we also use to terrify small children. It is much less complicated than it looks. The first and second lines are the signal and background parts respectively: the underlying time-distribution of the signal is exponential, while that of the background is modelled by a fraction ( $f_{\tau_{BG}}$ ) with lifetime (*e.g.* charm daughters), following  $e^{-t'/\tau_{BG}}$ , and a fraction without, following  $\delta(t')$ ;  $R$  is a double-Gaussian resolution function. The nine parameters shown are floated in the fit, so the background properties are fitted along with the signal: this is the reason for including the region  $3\sigma < |M^i - m_D| < 6\sigma$  in the fit, providing a background-rich sample which largely determines the background parameters.

(In fact there are eighteen fitted quantities, because there are two functions (2): one each for  $K^- \pi^+$  and  $K^- K^+$  decays. Our fit maximises the grand likelihood  $\mathcal{L} = \mathcal{L}_{K\pi} \cdot \mathcal{L}_{KK}$ , replacing the  $K^- K^+$  lifetime by  $\tau_{KK} = \tau_{K\pi} / (1 + y_{CP})$ . We find a null result, by the way:  $y_{CP} = (-0.5 \pm 1.0 \pm 0.8)\%$ .)

If you look closely at the resolution term, you'll notice that the proper-time error for each event is given by  $\sigma_t^i$ : an event-dependent quantity. Due to variations in track and D-vertex quality, the estimated vertexing and therefore proper-time errors vary from event to event, by a factor of a few. (Kinematic variations also play a role.) We scale these errors by global factors  $S$  (for the core Gaussian; close to 1) and  $S_{tail}$  (for the tail Gaussian;  $> 1$ ), but the full function  $R$  varies event-to-event, as does the signal fraction  $f_{SIG}^i$ . Any binned fit to the data would therefore need to have multidimensional bins—*many*-dimensional, for a complicated analysis like  $\sin 2\phi_1$ —and to avoid this, we resort to unbinned fits.

And so to the statistical issues raised by the timing-distribution fits:

1. *How do we estimate goodness-of-fit to our timing distributions?* While standard measures exist for binned fits, there is no accepted goodness-of-fit test for unbinned maximum likelihood. Some effort has recently gone into finding a method . . . if indeed it's possible [7]. In the meantime the lack of such a method is a nuisance, since we have nontrivial resolution functions which we fit from the data. How would we know if the functional form were wrong; and would it matter?
  - (a) In the case of  $y_{CP}$ , we perform extensive systematic checks by varying cuts, signal-to-background ratios and the like; and trace some biases to their origin by turning effects on and off in our detector Monte Carlo. This doesn't so much assure us that the fit is good, as that any variations of the fit, or problems with it, have a controlled effect on  $y_{CP}$ .
  - (b) For  $\sin 2\phi_1$ , we test the (*very* complicated) resolution function by using it in the measurement of B-lifetime: a simpler, and much-higher-statistics task, than the asymmetry fit. We check for biases by fitting null-asymmetry samples which are similar to our signal:  $B^+ \rightarrow J/\psi K^+$ , for example. And we compare our results to ensembles of fits to toy Monte Carlo datasets . . . although there may be less information in this last check than we once thought [7].

This is all good and necessary, and helps us (and our journal referees) to sleep at night. But if a decent goodness-of-fit test existed, we would obviously want to use that *as well*, and I for one would be glad if someone developed such a thing. Or could convince me not to worry about it.

2. *How should we combine our errors?* This is a problem we have largely postponed, as our unitarity angle analyses are still statistically limited. For  $y_{CP}$ , statistical and systematic errors are of equal magnitude, and we estimate the total error  $\Delta^2$  using the familiar  $\Delta^2 = \sigma^2(\text{stat.}) + \delta^2(\text{syst.})$ . Familiar is, of course, not the same as “correct”.

3. *How should we estimate confidence intervals?* For  $y_{CP}$ , what we *do* is to treat the  $\Delta$  just defined as a Gaussian-distributed error. (Yes, I know that’s an assumption.) For the unitarity angle analyses, confidence intervals . . . are a can of worms. See section 5

## 4.2 Searches for rare hadronic $B$ decays

Rare decay analyses are simpler, at least on the surface. These studies are motivated by the search for “direct” CP violation, *i.e.* CP asymmetry of decay amplitudes. Decays which are Cabibbo, CKM or colour-suppressed, or proceed via loop diagrams, are a good place to look for direct CPV: mechanisms with different CKM structure (such as Penguins and  $V_{ub}$  tree diagrams) can compete, with similar magnitudes; interference can lead to CP violation. For similar reasons, “New Physics” (non-Standard Model effects) can be expected to contribute, since the competing Standard Model processes are suppressed.

There are many (possible) rare decay modes, but the analyses tend to follow the pattern established in [8]. As an example, let’s take the somewhat simpler publication on  $B \rightarrow \eta' K$  and  $\eta' \pi$  decays [9].

### 4.2.1 What we do: the $B \rightarrow \eta' K$ ( $\eta' \pi$ ) analysis

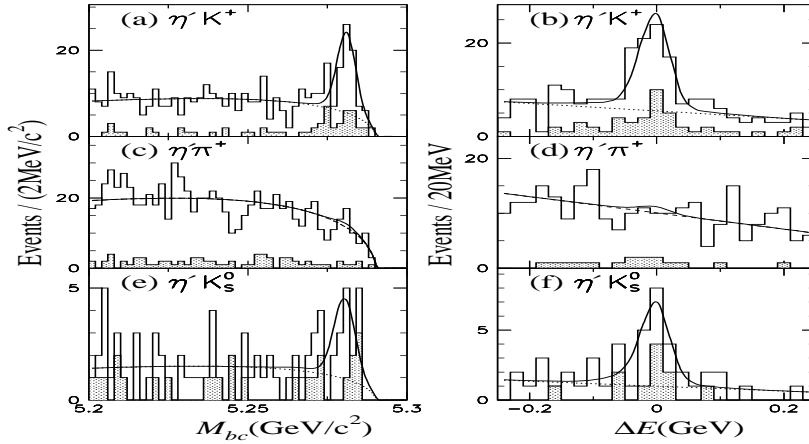


Fig. 2: From the  $\eta' h$  paper [9]:  $M_{bc}$  and  $\Delta E$  projections for  $\eta' K^+$ ,  $\eta' \pi^+$  and  $\eta' K_s^0$  events. The shaded histograms show  $\eta' \rightarrow \eta \pi^+ \pi^-$  events while the open histograms show  $\eta' \rightarrow \eta \pi^+ \pi^-$  and  $\eta' \rightarrow \rho \gamma$  combined. The superimposed curves show the fits to  $M_{bc}$  and  $\Delta E$  (solid) and the background component in the fit (dashed). Note the negligible  $\eta' \pi^+$  yield.

After making an event selection, continuum (*i.e.*  $e^+ e^- \rightarrow q \bar{q}$ ) backgrounds are suppressed using a likelihood ratio formed from (i) a Fisher discriminant of event shape variables, (ii) the production angle of the B-candidate, and (iii) for  $\eta' \rightarrow \rho \gamma$  decays, a helicity variable. A cut on the likelihood ratio gets rid of 70–90% of the background while keeping most of the signal. The signal is then isolated in  $M_{bc}$  and  $\Delta E$  (beam-constrained mass and energy-difference), which exploit the constrained kinematics of  $e^+ e^- \rightarrow \Upsilon(4S) \rightarrow B \bar{B}$  events: the results are shown in Fig. 2 for  $B^+ \rightarrow \eta' K^+$ ,  $\eta' \pi^+$  and  $B^0 \rightarrow \eta' K^0$ . Signal events appear in Gaussian peaks at  $M_{bc} = 5.28 \text{ GeV}/c^2$  and  $\Delta E = 0$ ; continuum backgrounds follow a phase-space-like function due to ARGUS [10] in  $M_{bc}$ , and a linear form in  $\Delta E$ . (Background shape parameters are set from appropriate sideband data, and cross-checked in the Monte Carlo simulation.)

Note that  $\eta' K^0$  is near the edge of our sensitivity, and that we see no  $\eta' \pi^+$  peak. We assess the significance of a yield using  $\Sigma \equiv \sqrt{-2 \ln(\mathcal{L}_0/\mathcal{L}_{max})}$ , where  $\mathcal{L}_{max}$  is the maximum likelihood returned by the fit (here, an unbinned fit to  $(M_{bc}, \Delta E)$ ), and  $\mathcal{L}_0$  is the likelihood at zero yield. For the cases  $\Sigma \geq 3 \equiv$  “**significant**”: we quote a *central value*, but no confidence interval (*e.g.*  $B^0 \rightarrow \eta' K^+$ );  $\Sigma < 3 \equiv$  “**not significant**”: we quote an *upper limit*, but no central value (*e.g.*  $B^0 \rightarrow \eta' \pi^+$ ).

Table 1: From the  $\eta' h$  paper [9]: Branching fraction ( $BF$ ) or 90% C.L. limit, and significance ( $\Sigma$ ) for Belle, compared with CLEO [11] and BABAR [12] results, and theoretical expectations [13, 14]. The branching fractions are in units of  $10^{-6}$ .

| Mode                          | This measurement( $BF$ ) | $\Sigma$ | CLEO                   | BABAR             | Theory |
|-------------------------------|--------------------------|----------|------------------------|-------------------|--------|
| $B^+ \rightarrow \eta' K^+$   | $79_{-11}^{+12} \pm 9$   | 12.0     | $80_{-9}^{+10} \pm 7$  | $62 \pm 18 \pm 8$ | 21–53  |
| $B^+ \rightarrow \eta' \pi^+$ | $< 7$                    | 0.0      | $< 12$                 | -                 | 1–3    |
| $B^0 \rightarrow \eta' K^0$   | $55_{-16}^{+19} \pm 8$   | 5.4      | $89_{-16}^{+18} \pm 9$ | $< 112$           | 20–50  |

There are some full-reporting issues here, but let’s set them aside. The *upper limit* (usually at 90% C.L.) is calculated using the notorious method of “*integrating the likelihood function*”, followed by addition of one unit of systematic error! (The results are shown in Table 1.) Where we measure CP-asymmetries between (say)  $B^+$  and  $B^-$  decays, intervals are constructed in the same spirit: for  $A_{CP}(B^\pm \rightarrow \eta' K^\pm) = x_{-\sigma^-}^{+\sigma^+}$  (stat.)  $_{-\delta^-}^{+\delta^+}$  (syst.), we set a 90% C.L. interval  $(x - 1.64\sigma^- - \delta^-, x + 1.64\sigma^+ + \delta^+)$ .

#### 4.22 Why what we do is not so bad

Like many of you, I have little good to say about this technique. A likelihood  $\mathcal{L}(\mu; x)$  for parameter(s)  $\mu$  given measurement(s)  $x$  is nothing other than the *probability density*  $p(x; \mu)$  to obtain the observed data, if the underlying parameter really were  $\mu$ . It is thus a density *in*  $x$ , and to “integrate” a density in the wrong variable is confused. (Try this with a Gaussian  $p(x; \mu) = (2\pi\sigma^2)^{-1/2} \exp(-(x - \mu)/2\sigma^2)$ , integrating  $\int d\mu p(x; \mu)$  for fixed  $x$ —as opposed to integrating it over  $x$ —and you’ll see what I mean.) And do not get me started on the addition of one unit of systematic error . . . at least until section 6.

Having got that out of the way, I now want to explain why this method is not as bad as it looks:

1. *It is easy and fast, and can be done with information already “in hand”.* All you need is the likelihood function. This is important, since the typical rare decay analysis needs to be published with some urgency: either because it is a first observation, or because some issue hangs on the measurement. ( $\eta' K(\pi)$  is of this kind: there is a theory/data discrepancy which might be a first hint of something exciting; see Table 1.) There is a kind of built-in obsolescence in these measurements too, due to continuing improvements in our luminosity: Physics Letters published our paper [9] on 4th October 2001: nine months after that date, we will have at least eight times its data on tape.
2. *For branching fraction measurements, it corresponds roughly to a Bayesian interval.* If we have some prior degree of belief  $p(\mu)$  concerning a parameter (here,  $\mu = B_{true}$ , the true branching fraction), then after the measurement  $x$  we may update this belief using Bayes’ Theorem

$$p(\mu|x) = p(x|\mu) \cdot p(\mu)/p(x), \quad (3)$$

where  $p(x|\mu) = \mathcal{L}(\mu; x)$  is the likelihood function, and  $p(x)$  may be recovered from the normalization. The posterior probability  $p(\mu|x)$ —our updated belief about  $\mu$ , following the measurement  $x$ —is a density in  $\mu$  by construction, and therefore *can* be integrated on  $\mu$ ; and “integrating the likelihood function” is equivalent to integrating (3) if the prior probability  $p(\mu)$  is constant.

Now a constant function is hard to defend as a serious prior, but there is a more subtle problem with this approach, to do with the special point  $\mu = 0$ . If our prior  $p(\mu)$  is truly constant, this means we are committed in advance to the belief that the branching fraction *does not vanish*, since  $S = \{0\}$  is a set of measure zero:  $\int_S d\mu p(\mu) = 0$  for any finite  $p(\mu)$ . If we were deriving a proper Bayesian credible interval for  $\mu$  we might well assign a *delta function* to the origin, allowing a (say) 10% belief that the decay is forbidden; the posterior for  $\mu = 0$  would then be nonzero.<sup>4</sup> The upper limit calculated *via*

$$\left[ \int_0^B d\mu p(x|\mu) \cdot p(\mu)/p(x) \right] / \left[ \int_0^\infty d\mu p(x|\mu) \cdot p(\mu)/p(x) \right] = 0.9 \quad (4)$$

<sup>4</sup>For  $P(0|x) = \int_S d\mu p(\mu|x) \ll 0.1$ , the point  $\mu = 0$  might lie *outside* the 90% interval, although naively integrating  $p(x|\mu)$  would never tell you that: (4) always yields an upper limit. That is, (4) is not a unified method for interval construction.

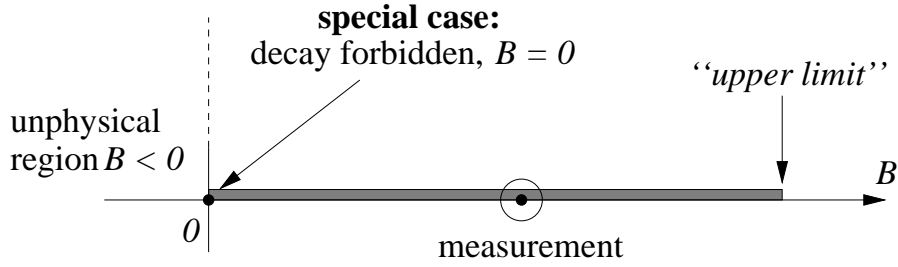


Fig. 3: A branching fraction measurement. The region  $B < 0$  is unphysical;  $B = 0$  is a special case. For a measurement  $x = B_{meas}$ , “integration of the likelihood function” (cf. (4)) over the shaded region yields an upper limit as shown.

would fall, but this is just a special case of the dependence of the limit on the prior. The special point  $\mu = 0$  is benign because it coincides with the physical boundary (Fig. 3) and is always on the *edge* of an interval, if it belongs to it at all. In section 5, we will see an example where this is not the case.

Upper limits derived in this way might differ by a factor of  $\approx 2$  from those we would obtain by a more rigorous (Bayesian) technique. Where there is already a tradition of quoting such limits, so that everyone “knows” what they mean (just as we know that “3 sigma” and “5 sigma” do not *really* mean 99.7% and 99.994% confidence), it would be hard to justify declaring war on the method. And the frequentist alternative is messy: it would require the construction of toy Monte Carlos for each and every decay mode and analysis (all of them differ subtly), to determine coverage . . . and we would probably find limits (again) within a factor 2, at the price of making lots of work (in parallel!) for lots of students.

### 4.3 Fits over a “flattish” background

These analyses are (to my mind) more straightforward: they involve fitting a lineshape over a smooth background, or at worst, interpreting the result of a background subtraction. Our analysis of prompt charmonium production [15] is an example. Selecting  $J/\psi$  events with centre-of-mass momentum  $p_\psi^* > 2.0 \text{ GeV}/c$ , above the kinematic limit for  $B \rightarrow J/\psi X$ , we measure the yield in the main “on-resonance” data, and in the smaller “off-resonance” sample, where  $\sqrt{s}$  is just below  $B\bar{B}$  threshold: too low to produce an  $\Upsilon(4S)$  meson. After scaling, correction and cross-checks, we subtract the yields to find the net number of  $\Upsilon(4S) \rightarrow J/\psi X$  decays to be  $-37 \pm 156$ : *i.e.* consistent with vanishing  $\mathcal{B}(\Upsilon(4S) \rightarrow J/\psi X)$ .

The error is dominated by the uncertainty on the off-resonance yield. We assume that this—and the full error—is distributed as a Gaussian, and use [16, table X] to determine the upper limit. This allows the negative yield to be treated in a rigorous way. Some approximation is involved by assuming Gaussian behaviour, and in principle we could model the subtraction in a toy Monte Carlo, and determine the limit using likelihood-ratio ordering from first principles. In my view the utility of spending 30 seconds looking up a table, and referring readers to an accessible paper, outweighs issues of principle here.

Having set an upper limit for  $\mathcal{B}(\Upsilon(4S) \rightarrow J/\psi X)$ , we assume that the observed prompt  $J/\psi$  are produced directly from  $e^+e^-$  annihilation. We look for prompt production of other charmonia, and find a significant yield for  $\psi(2S)$ , but not  $\chi_{c1,c2}$ : see Fig. 4. For  $\chi_{c1,c2}$  we set upper limits on  $\sigma(e^+e^- \rightarrow \chi_{c1,c2} X)$  using the Feldman-Cousins tables for Gaussians: the assumption of Gaussian errors is clearly reasonable. The same technique is used when treating small yields on a background in [17].

The search for (factorization-forbidden)  $B \rightarrow \chi_{c2} X$  decays [17] is unusual for Belle: a B-decay analysis where we fit a peak over a large, smooth background, so statistical questions are unproblematic. (Because of the inclusive nature of the decay, we cannot use the  $M_{bc}, \Delta E$  variables to wipe away the background.) There is however an interesting question concerning systematic errors. Because of the complicated (and overlapping) lineshapes used to fit  $\chi_{c1,c2} \rightarrow \gamma J/\psi$ , and the relatively large  $\chi_{c1}$  yield, the systematic error on the  $\chi_{c2}$  yield is substantial:  $\mathcal{B}(B \rightarrow \chi_{c2} X) = (1.80_{-0.28}^{+0.23} \pm 0.26) \times 10^{-3}$ .

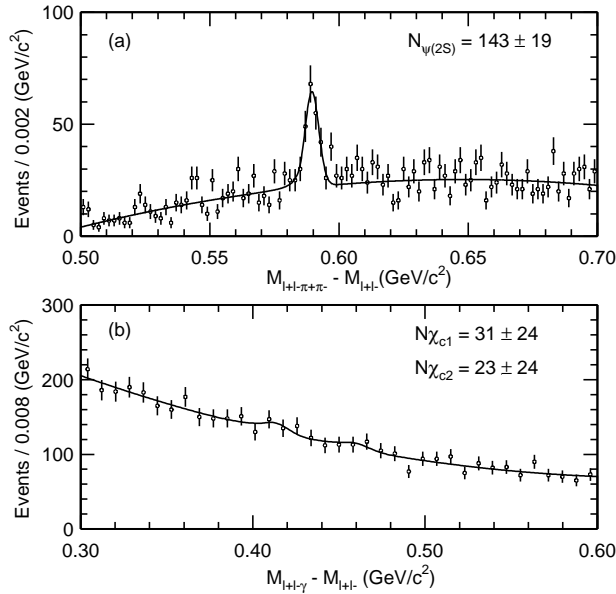


Fig. 4: From the prompt charmonium paper [15]: fits for the yield of  $\psi(2S) \rightarrow J/\psi \pi^+ \pi^-$  (upper plot) and  $\chi_{c1,2} \rightarrow \gamma J/\psi$  (lower plot) for  $p^* > 2.0 \text{ GeV}/c$ , above the limit for  $B \rightarrow (c\bar{c})_{res}$  decays.

The systematic error  $\pm 0.26$  is dominated by the choice of the fit function: we estimate the associated uncertainty using a large sample of reasonable (and some *un*-reasonable) variations to the fitting model. How should we “combine” the statistical and systematic errors in this case? I, for one, don’t know.

We *can* answer the following, more sharply posed question: is the yield “significant” even when the systematics are taken into account? (Phys. Rev. Lett. insists on “ $5\sigma$ ” significance before you are allowed to call something an “observation”.) All attempted variations to the fit gave yields with significance  $> 5\sigma$ , and we take this to be the relevant test: the statement that “the yield is inconsistent with fluctuations of the background” does not depend on some accidental feature of the fit, but is robust.

#### 4.4 Systematic-dominated analyses

When on the other hand systematic errors are dominant, we do not quote intervals ... and questions of “significance” tend not to arise. An example is our measurement of  $B \rightarrow X_s \gamma$  decays [18], where we find the underlying quark transition—the theoretically interesting process—to have a branching  $\mathcal{B}(b \rightarrow s\gamma) = (3.36 \pm 0.53 \text{ (stat.)} \pm 0.42 \text{ (syst.)} \pm_{-0.54}^{+0.50} \text{ (th.)}) \times 10^{-4}$ . There are plenty of issues of interpretation in such cases (they are *analyses* properly-so-called) but they are beyond the scope of this review.

### 5 THE NEW $B^0/\bar{B}^0 \rightarrow \pi^+ \pi^-$ ANALYSIS

As a final example, let’s consider a very difficult problem: the interpretation of our new  $B^0/\bar{B}^0 \rightarrow \pi^+ \pi^-$  result [5]. Like the  $\sin 2\phi_1$  analysis, this is a measurement of a time-dependent CP-violating asymmetry, with two differences:  $\pi^+ \pi^-$  is sensitive to the unitarity angle  $\phi_2$ ; and there may be direct CP violation in the decay. We fit

$$R_q(\Delta t) = \frac{e^{-|\Delta t|/\tau_B}}{4\tau_B} [1 + q \cdot \{A_{\pi\pi} \cos(\Delta m \Delta t) + S_{\pi\pi} \sin(\Delta m \Delta t)\}] \quad (5)$$



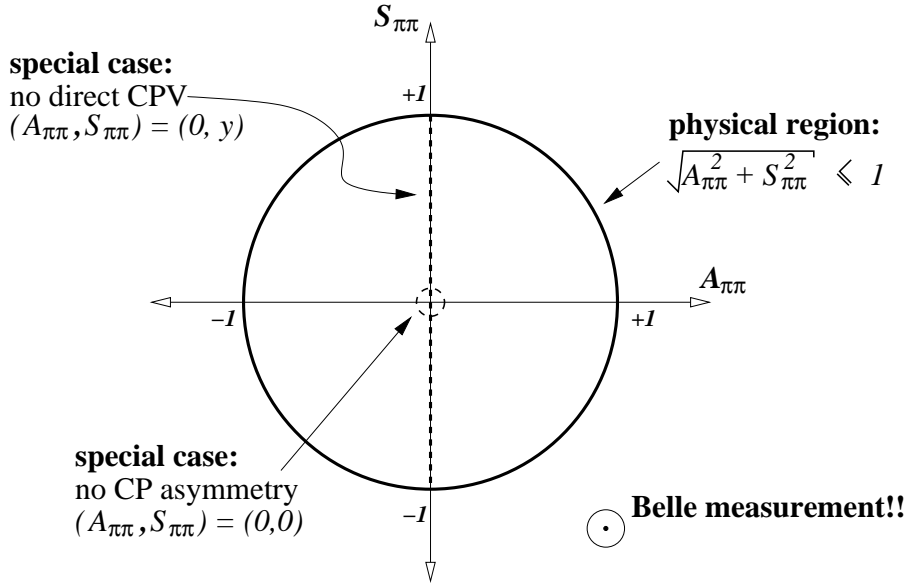


Fig. 5: The parameter space for the  $B^0/\bar{B}^0 \rightarrow \pi^+\pi^-$  asymmetry analysis, with the Belle measurement.

to the decay-time distribution for  $B^0$  ( $q = 1$ ) and  $\bar{B}^0$  ( $q = -1$ ), where  $\tau_B$  and  $\Delta m$  are the lifetime and mass-splitting of the eigenstates  $|B_\pm\rangle = p|B^0\rangle \pm q|\bar{B}^0\rangle$ . The coefficients are given by

$$A_{\pi\pi} = \frac{|\lambda|^2 - 1}{|\lambda|^2 + 1}, \quad S_{\pi\pi} = \frac{2\Im\lambda}{|\lambda|^2 + 1}, \quad \lambda = \frac{q}{p} \cdot \frac{A(\bar{B}^0 \rightarrow \pi^+\pi^-)}{A(B^0 \rightarrow \pi^+\pi^-)} \quad (6)$$

A value  $A_{\pi\pi} \neq 0$  corresponds to direct CP violation in the decay.

The parameter space, Fig. 5, is wonderfully complicated. There are two kinds of special region: the null asymmetry point  $(A_{\pi\pi}, S_{\pi\pi}) = (0, 0)$ , and the *line*  $(A_{\pi\pi}, S_{\pi\pi}) = (0, y)$ , where there is no direct CPV ( $A_{\pi\pi} = 0$ ). The physical boundary  $\sqrt{A_{\pi\pi}^2 + S_{\pi\pi}^2} = 1$  forms a ring around these regions ... and the Belle measurement  $A_{\pi\pi} = +0.94_{-0.31}^{+0.25} \pm 0.09$ ,  $S_{\pi\pi} = -1.21_{-0.31}^{+0.25} {}_{-0.13}^{+0.16}$  lies *outside it!*<sup>5</sup> If we want to determine a confidence interval in the true parameters  $(a_{\pi\pi}, s_{\pi\pi})$ , we immediately run into difficulty:

1. *Frequentist method:* If we study the fit using Monte Carlo events, we find that the fitted error on  $(A_{\pi\pi}, S_{\pi\pi})$  *varies* from virtual experiment to virtual experiment, by a factor of a few. So in constructing a toy Monte Carlo to model the fit—and determine confidence intervals—how do we generate the observed  $(A_{\pi\pi}, S_{\pi\pi})$ , given some underlying parameters  $(a_{\pi\pi}, s_{\pi\pi})$ ?
  - (a) If we use the *measured* errors, what if we “get lucky” and they are unusually small? Can we trust a result that says our value is “inconsistent” with fluctuations from  $(0, 0)$ ?
  - (b) If we use the *distribution* of errors from the full Monte Carlo, then the *actual* errors returned by the fit are never used in the analysis. It seems paradoxical to “throw away” a fitted error.
2. *Bayesian method:* Here the community expectation is less definite than for  $B^0/\bar{B}^0 \rightarrow J/\psi K_S^0$ , so any explicit prior would be controversial. But because of the configuration of the physically interesting points and the physical boundary, an explicit prior seems unavoidable. For suppose we tried to form a pseudo-Bayesian interval by “integrating the likelihood function” on the physical region, expanding outwards from the measured value (Fig. 6). We want to know if we reach the point of null asymmetry  $(0, 0)$  before or after we have exhausted (say) 99.7% of the integral: if the

<sup>5</sup>Since the experimental quantity is an asymmetry, this is not possible for pure signal; *cf.* a  $\beta$ -decay endpoint analysis, where resolution effects can give  $m_\nu^2 < 0$ . In the presence of background events, values outside the physical boundary can occur.

99.7% interval does not include  $(0, 0)$ , surely we can say that we exclude it at some level . . . and write off the improper treatment of the prior as a minor detail (*cf.* section 4.22)?

This is not possible in general, because of the special region  $A_{\pi\pi} = 0$ . If we allow some prior probability for indirect CPV ( $S_{\pi\pi} \neq 0$ ) in the absence of direct CP violation, there is a “delta-line function”  $\delta(a) \cdot f(s)$  running along the  $S_{\pi\pi}$  axis. In general a credible interval will intersect this function *before reaching*  $(0, 0)$ . Thus there is no way to determine the probability content of the interval without making an explicit commitment as to the prior  $\delta(a) \cdot f(s)$ . (A “flat” prior means  $f(s) = 0$ .) I can see no way around this problem: any prior must be explicit.

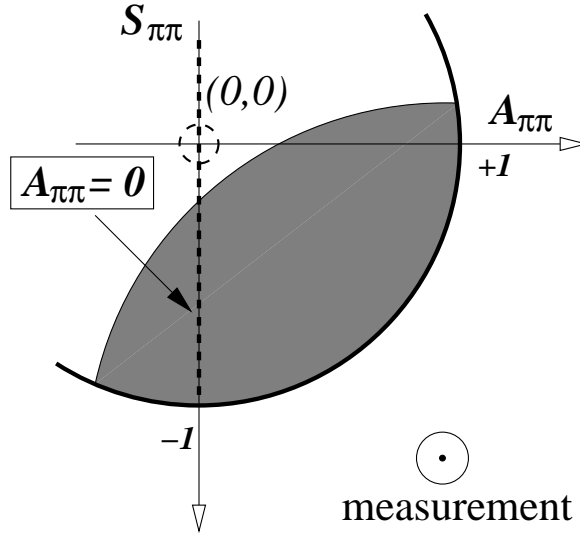


Fig. 6: Detail of Fig. 5, showing an interval (shaded) constructed with a “flat prior” on the physical region. Note that the interval intersects the special region  $A_{\pi\pi} = 0$  before reaching the null asymmetry point  $(0, 0)$ .

## 6 SOME OPEN QUESTIONS

### 6.1 How do we quote rare decay results?

In section 4.2, in all the excitement about priors, I almost forgot to notice a glaring inconsistency. We require significance  $\Sigma > 3$  before we will quote a central value, corresponding to a 99.7% confidence requirement . . . but we quote 90% C.L. upper limits! If a value falls between these two levels, we should not be able to say anything at all: in fact we are let off the hook because our interval-setting is not unified. (As noted above, the integral method (4) always gives an upper limit.) In defence of Belle, we are following community practice and expectations here, and those expectations are incoherent. My provisional idea for a way around the problem is to

- *always* quote the central value;
- construct 99.7% C.L. intervals in a unified manner (*i.e.* going over continuously from upper limits to central intervals, the way the so-called Feldman-Cousins intervals [16] do);
- use these intervals in place of the  $\Sigma > 3$  test;
- quote 90% C.L. intervals *as well*, because people expect them;
- if people query the use of 3 numbers/intervals, rather than 1, explain;
- if people *object* to the use of 3 numbers, resort to violence.

This would be consistent, but it is a utopian scheme . . . and I suspect it would take a dictator to implement it (*cf.* section 3). I will try to raise consciousness on this issue at Belle, but it really is a

community problem, and we should try to think up some way for all of us to get to “there” from “here”. The three-values approach I’ve suggested may not be the best way.

Note also that a unified treatment of “significance” and confidence intervals begs another question:

## 6.2 How do we combine statistical and systematic errors?

In sections 4.1 and 4.3 I noted cases where we are able to avoid this question, but this is not general. The technique used for the rare decays (section 4.2) is peculiar: the integration method is pseudo-Bayesian, as discussed at some length; it may not be quite so obvious, but the practice of inflating the confidence intervals by “one sigma” of the systematic error is pseudo-Frequentist. It treats the systematic error as a nuisance parameter with range  $[-\delta^-, +\delta^+]$ , and demands that our confidence interval provides coverage for all values in that region. To my mind this is almost exactly the wrong way around:

- I think our prior beliefs about branching fractions and CP violation are not of interest (or at least, do not belong in papers), which suggests a Frequentist approach; whereas
- our beliefs about our systematic errors surely are relevant—everyone knows they come down to a question of judgement—which suggests a Bayesian approach.

A complicating factor is that not all “systematic errors” are the same kind of thing:

- Particle ID efficiencies are measured on control samples, which have statistical uncertainties: it seems reasonable to use some averaging method when assessing the effect on any final interval.
- The unknown phases of resonances on Dalitz plots (for example) are at the other extreme: we have no business making assumptions about them, and if they significantly affect a result we should require any confidence interval to provide coverage for all possible values.
- The choice of parameterization for a fit function is unlike both of the previous examples.

Needless to say, these musings do not constitute a policy, much less a recipe. I suspect that intervals of the form  $(x - 1.64\sigma^- - \delta^-, x + 1.64\sigma^+ + \delta^+)$  will remain with us for some time.

## 6.3 What should we do about the $\pi^+\pi^-$ analysis?

The consensus from the floor during this talk was that in a problem of this difficulty, rigour is essential: only the two extreme approaches make sense, *viz.*

1. a Frequentist calculation from first principles, and
2. a full, openly subjective Bayesian analysis.

Regarding the second of these, I am already convinced (see section 5). As for the Frequentist calculation, I am indebted to my colleagues at the conference for some ingenious suggestions on the proper treatment of the errors. I hope to experiment with them, together with my Belle colleagues, before the summer.

## 7 CONCLUSION

The Belle collaboration performs a large range of analyses, using a range of statistical approaches. Some of the methods are open to question, although there is a clear trade-off between utility and statistical rigour in some cases. The expectations of the broader physics community also play a role. In the case of the  $B^0/\bar{B}^0 \rightarrow \pi^+\pi^-$  analysis, the statistical environment is unusually difficult, and rigorous (rather than approximate) methods are required.

## ACKNOWLEDGEMENTS

I would like to thank the conference organizers for arranging an instructive and entertaining programme for this meeting. In particular, I am indebted to the conference secretary Linda Wilkinson, for her help in

recovering lost slides and other material. To my colleagues on Belle, and the collaboration management: thanks for their patience with my questions and criticisms on statistical matters, and for the freedom I am habitually given to speak my mind in public.

## References

- [1] M. Kobayashi and T. Maskawa, *Prog. Theor. Phys.* **49**, 652–657 (1973).
- [2] K. Abe *et al.*, *Phys. Rev. Lett.* **87**, 091802 (2001).
- [3] B. Aubert *et al.*, *Phys. Rev. Lett.* **87**, 091801 (2001).
- [4] A. Abashian *et al.*, *Nucl. Instr. Meth. A* **479**, 117–232 (2002).
- [5] K. Abe *et al.*, hep-ex/0204002, accepted for publication in *Phys. Rev. Lett.*
- [6] K. Abe *et al.*, *Phys. Rev. Lett.* **88**, 162001 (2002).
- [7] Kay Kinoshita, in these proceedings.
- [8] K. Abe *et al.*, *Phys. Rev. Lett.* **87**, 101801 (2001).
- [9] K. Abe *et al.*, *Phys. Lett. B* **517**, 309–318 (2001).
- [10] H. Albrecht *et al.*, *Phys. Lett. B* **241**, 278–282 (1990).
- [11] S. J. Richichi *et al.*, *Phys. Rev. Lett.* **85** (2000) 520–524.
- [12] T. Champion, in *Proc. XXXth Int. Conf. on High Energy Phys. (ICHEP2000)*, edited by C.S. Lim and T. Yamanaka, World Scientific, Singapore (2001).
- [13] A. Ali, G. Kramer, and C.-D. Lu, *Phys. Rev. D* **58** (1998) 094009; Y.-H. Chen, H.-Y. Cheng, B. Tseng, and K.-C. Yang, *Phys. Rev. D* **60** (1999) 094014; H.-Y. Cheng and K.C. Yang, *Phys. Rev. D* **62** (2000) 054029.
- [14] E. Kou and A. I. Sanda, *Phys. Lett. B* **525**, 240–248 (2002).
- [15] K. Abe *et al.*, *Phys. Rev. Lett.* **88**, 052001 (2002).
- [16] G. J. Feldman and R. D. Cousins, *Phys. Rev. D* **57**, 3873–3889 (1998).
- [17] K. Abe *et al.*, *Phys. Rev. Lett.* **89**, 011803 (2002).
- [18] K. Abe *et al.*, *Phys. Lett. B* **511**, 151–158 (2001).