

# Does Psi Exist? Replicable Evidence for an Anomalous Process of Information Transfer

Daryl J. Bem and Charles Honorton

Most academic psychologists do not yet accept the existence of psi, anomalous processes of information or energy transfer (such as telepathy or other forms of extrasensory perception) that are currently unexplained in terms of known physical or biological mechanisms. We believe that the replication rates and effect sizes achieved by one particular experimental method, the ganzfeld procedure, are now sufficient to warrant bringing this body of data to the attention of the wider psychological community. Competing meta-analyses of the ganzfeld database are reviewed, 1 by R. Hyman (1985), a skeptical critic of psi research, and the other by C. Honorton (1985), a parapsychologist and major contributor to the ganzfeld database. Next the results of 11 new ganzfeld studies that comply with guidelines jointly authored by R. Hyman and C. Honorton (1986) are summarized. Finally, issues of replication and theoretical explanation are discussed.

The term *psi* denotes anomalous processes of information or energy transfer, processes such as telepathy or other forms of extrasensory perception that are currently unexplained in terms of known physical or biological mechanisms. The term is purely descriptive: It neither implies that such anomalous phenomena are paranormal nor connotes anything about their underlying mechanisms.

Does psi exist? Most academic psychologists don't think so. A survey of more than 1,100 college professors in the United States found that 55% of natural scientists, 66% of social scientists (excluding psychologists), and 77% of academics in the arts, humanities, and education believed that ESP is either an established fact or a likely possibility. The comparable figure for psy-

chologists was only 34%. Moreover, an equal number of psychologists declared ESP to be an impossibility, a view expressed by only 2% of all other respondents (Wagner & Monnet, 1979).

We psychologists are probably more skeptical about psi for several reasons. First, we believe that extraordinary claims require extraordinary proof. And although our colleagues from other disciplines would probably agree with this dictum, we are more likely to be familiar with the methodological and statistical requirements for sustaining such claims, as well as with previous claims that failed either to meet those requirements or to survive the test of successful replication. Even for ordinary claims, our conventional statistical criteria are conservative. The sacred  $p = .05$  threshold is a constant reminder that it is far more sinful to assert that an effect exists when it does not (the Type I error) than to assert that an effect does not exist when it does (the Type II error).

Second, most of us distinguish sharply between phenomena whose explanations are merely obscure or controversial (e.g., hypnosis) and phenomena such as psi that appear to fall outside our current explanatory framework altogether. (Some would characterize this as the difference between the unexplained and the inexplicable.) In contrast, many laypersons treat all exotic psychological phenomena as epistemologically equivalent; many even consider *déjà vu* to be a psychic phenomenon. The blurring of this critical distinction is aided and abetted by the mass media, "new age" books and mind-power courses, and "psychic" entertainers who present both genuine hypnosis and fake "mind reading" in the course of a single performance. Accordingly, most laypersons would not have to revise their conceptual model of reality as radically as we would in order to assimilate the existence of psi. For us, psi is simply more extraordinary.

Finally, research in cognitive and social psychology has sensitized us to the errors and biases that plague intuitive attempts to draw valid inferences from the data of everyday experience (Gilovich, 1991; Nisbett & Ross, 1980; Tversky & Kahneman, 1971). This leads us to give virtually no probative weight to anecdotal or journalistic reports of psi, the main source cited by

Daryl J. Bem, Department of Psychology, Cornell University; Charles Honorton, Department of Psychology, University of Edinburgh, Edinburgh, Scotland.

Sadly, Charles Honorton died of a heart attack on November 4, 1992, 9 days before this article was accepted for publication. He was 46. Parapsychology has lost one of its most valued contributors. I have lost a valued friend.

This collaboration had its origins in a 1983 visit I made to Honorton's Psychophysical Research Laboratories (PRL) in Princeton, New Jersey, as one of several outside consultants brought in to examine the design and implementation of the experimental protocols.

Preparation of this article was supported, in part, by grants to Charles Honorton from the American Society for Psychical Research and the Parapsychology Foundation, both of New York City. The work at PRL summarized in the second half of this article was supported by the James S. McDonnell Foundation of St. Louis, Missouri, and by the John E. Fetzer Foundation of Kalamazoo, Michigan.

Helpful comments on drafts of this article were received from Deborah Delano, Edwin May, Donald McCarthy, Robert Morris, John Palmer, Robert Rosenhal, Lee Ross, Jessica Uris, Philip Zimbardo, and two anonymous reviewers.

Correspondence concerning this article should be addressed to Daryl J. Bem, Department of Psychology, Uris Hall, Cornell University, Ithaca, New York 14853. Electronic mail may be sent to d.bem@cornell.edu.

our academic colleagues as evidence for their beliefs about psi (Wagner & Monnet, 1979).

Ironically, however, psychologists are probably not more familiar than others with recent experimental research on psi. Like most psychological research, parapsychological research is reported primarily in specialized journals; unlike most psychological research, however, contemporary parapsychological research is not usually reviewed or summarized in psychology's textbooks, handbooks, or mainstream journals. For example, only 1 of 64 introductory psychology textbooks recently surveyed even mentions the experimental procedure reviewed in this article, a procedure that has been in widespread use since the early 1970s (Roig, Icochea, & Curzucoli, 1991). Other secondary sources for nonspecialists are frequently inaccurate in their descriptions of parapsychological research. (For discussions of this problem, see Child, 1985, and Palmer, Honorton, & Utts, 1989.)

This situation may be changing. Discussions of modern psi research have recently appeared in a widely used introductory textbook (Atkinson, Atkinson, Smith, & Bem, 1990, 1993), two mainstream psychology journals (Child, 1985; Rao & Palmer, 1987), and a scholarly but accessible book for nonspecialists (Broughton, 1991). The purpose of the present article is to supplement these broader treatments with a more detailed, meta-analytic presentation of evidence issuing from a single experimental method: the *ganzfeld* procedure. We believe that the replication rates and effect sizes achieved with this procedure are now sufficient to warrant bringing this body of data to the attention of the wider psychological community.

### The Ganzfeld Procedure

By the 1960s, a number of parapsychologists had become dissatisfied with the familiar ESP testing methods pioneered by J. B. Rhine at Duke University in the 1930s. In particular, they believed that the repetitive forced-choice procedure in which a subject repeatedly attempts to select the correct "target" symbol from a set of fixed alternatives failed to capture the circumstances that characterize reported instances of psi in everyday life.

Historically, psi has often been associated with meditation, hypnosis, dreaming, and other naturally occurring or deliberately induced altered states of consciousness. For example, the view that psi phenomena can occur during meditation is expressed in most classical texts on meditative techniques; the belief that hypnosis is a psi-conductive state dates all the way back to the days of early mesmerism (Dingwall, 1968); and cross-cultural surveys indicate that most reported "real-life" psi experiences are mediated through dreams (Green, 1960; Prasad & Stevenson, 1968; L. E. Rhine, 1962; Sannwald, 1959).

There are now reports of experimental evidence consistent with these anecdotal observations. For example, several laboratory investigators have reported that meditation facilitates psi performance (Honorton, 1977). A meta-analysis of 25 experiments on hypnosis and psi conducted between 1945 and 1981 in 10 different laboratories suggests that hypnotic induction may also facilitate psi performance (Schechter, 1984). And dream-mediated psi was reported in a series of experiments conducted at Maimonides Medical Center in New York and

published between 1966 and 1972 (Child, 1985; Ullman, Krippner, & Vaughan, 1973).

In the Maimonides dream studies, two subjects—a "receiver" and a "sender"—spent the night in a sleep laboratory. The receiver's brain waves and eye movements were monitored as he or she slept in an isolated room. When the receiver entered a period of REM sleep, the experimenter pressed a buzzer that signaled the sender—under the supervision of a second experimenter—to begin a sending period. The sender would then concentrate on a randomly chosen picture (the "target") with the goal of influencing the content of the receiver's dream.

Toward the end of the REM period, the receiver was awakened and asked to describe any dream just experienced. This procedure was repeated throughout the night with the same target. A transcription of the receiver's dream reports was given to outside judges who blindly rated the similarity of the night's dreams to several pictures, including the target. In some studies, similarity ratings were also obtained from the receivers themselves. Across several variations of the procedure, dreams were judged to be significantly more similar to the target pictures than to the control pictures in the judging sets (failures to replicate the Maimonides results were also reviewed by Child, 1985).

These several lines of evidence suggested a working model of psi in which psi-mediated information is conceptualized as a weak signal that is normally masked by internal somatic and external sensory "noise." By reducing ordinary sensory input, these diverse psi-conductive states are presumed to raise the signal-to-noise ratio, thereby enhancing a person's ability to detect the psi-mediated information (Honorton, 1969, 1977). To test the hypothesis that a reduction of sensory input itself facilitates psi performance, investigators turned to the ganzfeld procedure (Braud, Wood, & Braud, 1975; Honorton & Harper, 1974; Parker, 1975), a procedure originally introduced into experimental psychology during the 1930s to test propositions derived from gestalt theory (Avant, 1965; Metzger, 1930).

Like the dream studies, the psi ganzfeld procedure has most often been used to test for telepathic communication between a sender and a receiver. The receiver is placed in a reclining chair in an acoustically isolated room. Translucent ping-pong ball halves are taped over the eyes and headphones are placed over the ears; a red floodlight directed toward the eyes produces an undifferentiated visual field, and white noise played through the headphones produces an analogous auditory field. It is this homogeneous perceptual environment that is called the *Ganzfeld* ("total field"). To reduce internal somatic "noise," the receiver typically also undergoes a series of progressive relaxation exercises at the beginning of the ganzfeld period.

The sender is sequestered in a separate acoustically isolated room, and a visual stimulus (art print, photograph, or brief videotaped sequence) is randomly selected from a large pool of such stimuli to serve as the target for the session. While the sender concentrates on the target, the receiver provides a continuous verbal report of his or her ongoing imagery and mentation, usually for about 30 minutes. At the completion of the ganzfeld period, the receiver is presented with several stimuli (usually four) and, without knowing which stimulus was the target, is asked to rate the degree to which each matches the imagery and mentation experienced during the ganzfeld period. If the receiver assigns the highest rating to the target stimulus, it

is scored as a "hit." Thus, if the experiment uses judging sets containing four stimuli (the target and three decoys or control stimuli), the hit rate expected by chance is .25. The ratings can also be analyzed in other ways; for example, they can be converted to ranks or standardized scores within each set and analyzed parametrically across sessions. And, as with the dream studies, the similarity ratings can also be made by outside judges using transcripts of the receiver's mentation report.

#### Meta-Analyses of the Ganzfeld Database

In 1985 and 1986, the *Journal of Parapsychology* devoted two entire issues to a critical examination of the ganzfeld database. The 1985 issue comprised two contributions: (a) a meta-analysis and critique by Ray Hyman (1985), a cognitive psychologist and skeptical critic of parapsychological research, and (b) a competing meta-analysis and rejoinder by Charles Honorton (1985), a parapsychologist and major contributor to the ganzfeld database. The 1986 issue contained four commentaries on the Hyman-Honorton exchange, a joint communique by Hyman and Honorton, and six additional commentaries on the joint communique itself. We summarize the major issues and conclusions here.

#### Replication Rates

*Rates by study.* Hyman's meta-analysis covered 42 psi ganzfeld studies reported in 34 separate reports written or published from 1974 through 1981. One of the first problems he discovered in the database was multiple analysis. As noted earlier, it is possible to calculate several indexes of psi performance in a ganzfeld experiment and, furthermore, to subject those indexes to several kinds of statistical treatment. Many investigators reported multiple indexes or applied multiple statistical tests without adjusting the criterion significance level for the number of tests conducted. Worse, some may have "shopped" among the alternatives until finding one that yielded a significantly successful outcome. Honorton agreed that this was a problem.

Accordingly, Honorton applied a uniform test on a common index across all studies from which the pertinent datum could be extracted, regardless of how the investigators had analyzed the data in the original reports. He selected the proportion of hits as the common index because it could be calculated for the largest subset of studies: 28 of the 42 studies. The hit rate is also a conservative index because it discards most of the rating information; a second place ranking—a "near miss"—receives no more credit than a last place ranking. Honorton then calculated the exact binomial probability and its associated  $z$  score for each study.

Of the 28 studies, 23 (82%) had positive  $z$  scores ( $p = 4.6 \times 10^{-4}$ , exact binomial test with  $p = q = .5$ ). Twelve of the studies (43%) had  $z$  scores that were independently significant at the 5% level ( $p = 3.5 \times 10^{-9}$ , binomial test with 28 studies,  $p = .05$ , and  $q = .95$ ), and 7 of the studies (25%) were independently significant at the 1% level ( $p = 9.8 \times 10^{-9}$ ). The composite Stouffer  $z$  score across the 28 studies was 6.60 ( $p = 2.1 \times 10^{-11}$ ).<sup>1</sup> A more conservative estimate of significance can be obtained by including 10 additional studies that also used the relevant judging procedure but did not report hit rates. If these studies

are assigned a mean  $z$  score of zero, the Stouffer  $z$  across all 38 studies becomes 5.67 ( $p = 7.3 \times 10^{-9}$ ).

Thus, whether one considers only the studies for which the relevant information is available or includes a null estimate for the additional studies for which the information is not available, the aggregate results cannot reasonably be attributed to chance. And, by design, the cumulative outcome reported here cannot be attributed to the inflation of significance levels through multiple analysis.

*Rates by laboratory.* One objection to estimates such as those just described is that studies from a common laboratory are not independent of one another (Parker, 1978). Thus, it is possible for one or two investigators to be disproportionately responsible for a high replication rate, whereas other, independent investigators are unable to obtain the effect.

The ganzfeld database is vulnerable to this possibility. The 28 studies providing hit rate information were conducted by investigators in 10 different laboratories. One laboratory contributed 9 of the studies, Honorton's own laboratory contributed 5, 2 other laboratories contributed 3 each, 2 contributed 2 each, and the remaining 4 laboratories each contributed 1. Thus, half of the studies were conducted by only 2 laboratories, 1 of them Honorton's own.

Accordingly, Honorton calculated a separate Stouffer  $z$  score for each laboratory. Significantly positive outcomes were reported by 6 of the 10 laboratories, and the combined  $z$  score across laboratories was 6.16 ( $p = 3.6 \times 10^{-10}$ ). Even if all of the studies conducted by the 2 most prolific laboratories are discarded from the analysis, the Stouffer  $z$  across the 8 other laboratories remains significant ( $z = 3.67$ ,  $p = 1.2 \times 10^{-4}$ ). Four of these studies are significant at the 1% level ( $p = 9.2 \times 10^{-6}$ , binomial test with 14 studies,  $p = .01$ , and  $q = .99$ ), and each was contributed by a different laboratory. Thus, even though the total number of laboratories in this database is small, most of them have reported significant studies, and the significance of the overall effect does not depend on just one or two of them.

#### Selective Reporting

In recent years, behavioral scientists have become increasingly aware of the "file-drawer" problem: the likelihood that successful studies are more likely to be published than unsuccessful studies, which are more likely to be consigned to the file drawers of their disappointed investigators (Bozarth & Roberts, 1972; Sterling, 1959). Parapsychologists were among the first to become sensitive to the problem, and, in 1975, the Parapsychological Association Council adopted a policy opposing the selective reporting of positive outcomes. As a consequence, negative findings have been routinely reported at the association's meetings and in its affiliated publications for almost two decades. As has already been shown, more than half of the ganzfeld studies included in the meta-analysis yielded outcomes whose significance falls short of the conventional .05 level.

A variant of the selective reporting problem arises from what

<sup>1</sup> Stouffer's  $z$  is computed by dividing the sum of the  $z$  scores for the individual studies by the square root of the number of studies (Rosenthal, 1978).

Hyman (1985) has termed the "retrospective study." An investigator conducts a small set of exploratory trials. If they yield null results, they remain exploratory and never become part of the official record; if they yield positive results, they are defined as a study after the fact and are submitted for publication. In support of this possibility, Hyman noted that there are more significant studies in the database with fewer than 20 trials than one would expect under the assumption that, all other things being equal, statistical power should increase with the square root of the sample size. Although Honorton questioned the assumption that "all other things" are in fact equal across the studies and disagreed with Hyman's particular statistical analysis, he agreed that there is an apparent clustering of significant studies with fewer than 20 trials. (Of the complete ganzfeld database of 42 studies, 8 involved fewer than 20 trials, and 6 of those studies reported statistically significant results.)

Because it is impossible, by definition, to know how many unknown studies—exploratory or otherwise—are languishing in file drawers, the major tool for estimating the seriousness of selective reporting problems has become some variant of Rosenthal's file-drawer statistic, an estimate of how many unreported studies with  $z$  scores of zero would be required to exactly cancel out the significance of the known database (Rosenthal, 1979). For the 28 direct-hit ganzfeld studies alone, this estimate is 4.23 fugitive studies, a ratio of unreported-to-reported studies of approximately 1.5:1. When it is recalled that a single ganzfeld session takes over an hour to conduct, it is not surprising that—despite his concern with the retrospective study problem—Hyman concurred with Honorton and other participants in the published debate that selective reporting cannot plausibly account for the overall statistical significance of the psi ganzfeld database (Hyman & Honorton, 1986).<sup>2</sup>

### Methodological Flaws

If the most frequent criticism of parapsychology is that it has not produced a replicable psi effect, the second most frequent criticism is that many, if not most, psi experiments have inadequate controls and procedural safeguards. A frequent charge is that positive results emerge primarily from initial, poorly controlled studies and then vanish as better controls and safeguards are introduced.

Fortunately, meta-analysis provides a vehicle for empirically evaluating the extent to which methodological flaws may have contributed to artifactual positive outcomes across a set of studies. First, ratings are assigned to each study that index the degree to which particular methodological flaws are or are not present; these ratings are then correlated with the studies' outcomes. Large positive correlations constitute evidence that the observed effect may be artifactual.

In psi research, the most fatal flaws are those that might permit a subject to obtain the target information in normal sensory fashion, either inadvertently or through deliberate cheating. This is called the problem of *sensory leakage*. Another potentially serious flaw is inadequate randomization of target selection.

*Sensory leakage.* Because the ganzfeld is itself a perceptual isolation procedure, it goes a long way toward eliminating potential sensory leakage during the ganzfeld portion of the ses-

sion. There are, however, potential channels of sensory leakage after the ganzfeld period. For example, if the experimenter who interacts with the receiver knows the identity of the target, he or she could bias the receiver's similarity ratings in favor of correct identification. Only one study in the database contained this flaw, a study in which subjects actually performed slightly below chance expectation. Second, if the stimulus set given to the receiver for judging contains the actual physical target handled by the sender during the sending period, there might be cues (e.g., fingerprints, smudges, or temperature differences) that could differentiate the target from the decoys. Moreover, the process of transferring the stimulus materials to the receiver's room itself opens up other potential channels of sensory leakage. Although contemporary ganzfeld studies have eliminated both of these possibilities by using duplicate stimulus sets, some of the earlier studies did not.

Independent analyses by Hyman and Honorton agreed that there was no correlation between inadequacies of security against sensory leakage and study outcome. Honorton further reported that if studies that failed to use duplicate stimulus sets were discarded from the analysis, the remaining studies are still highly significant (Stouffer  $z = 4.35$ ,  $p = 6.8 \times 10^{-9}$ ).

*Randomization.* In many psi experiments, the issue of target randomization is critical because systematic patterns in inadequately randomized target sequences might be detected by subjects during a session or might match subjects' preexisting response biases. In a ganzfeld study, however, randomization is a much less critical issue because only one target is selected during the session and most subjects serve in only one session. The primary concern is simply that all the stimuli within each judging set be sampled uniformly over the course of the study. Similar considerations govern the second randomization, which takes place after the ganzfeld period and determines the sequence in which the target and decoys are presented to the receiver (or external judge) for judging.

Nevertheless, Hyman and Honorton disagreed over the findings here. Hyman claimed there was a correlation between flaws of randomization and study outcome; Honorton claimed there was not. The sources of this disagreement were in conflicting definitions of flaw categories, in the coding and assignment of flaw ratings to individual studies, and in the subsequent statistical treatment of those ratings.

Unfortunately, there have been no ratings of flaws by independent raters who were unaware of the studies' outcomes (Morris, 1991). Nevertheless, none of the contributors to the subsequent debate concurred with Hyman's conclusion, whereas four nonparapsychologists—two statisticians and two psychologists—explicitly concurred with Honorton's conclusion (Harris & Rosenthal, 1988b; Saunders, 1985; Uts, 1991a). For example, Harris and Rosenthal (one of the pioneers in the use of meta-analysis in psychology) used Hyman's own flaw ratings and failed to find any significant relationships between flaws and study outcomes in each of two separate analyses:

<sup>2</sup> A 1980 survey of parapsychologists uncovered only 19 completed but unreported ganzfeld studies. Seven of these had achieved significantly positive results, a proportion (.37) very similar to the proportion of independently significant studies in the meta-analysis (.43) (Blackmore, 1980).

"Our analysis of the effects of flaws on study outcome lends no support to the hypothesis that Ganzfeld research results are a significant function of the set of flaw variables" (1988b, p. 3; for a more recent exchange regarding Hyman's analysis, see Hyman, 1991; Uts, 1991a, 1991b).

### Effect Size

Some critics of parapsychology have argued that even if current laboratory-produced psi effects turn out to be replicable and nonartificial, they are too small to be of theoretical interest or practical importance. We do not believe this to be the case for the psi ganzfeld effect.

In psi ganzfeld studies, the hit rate itself provides a straightforward descriptive measure of effect size, but this measure cannot be compared directly across studies because they do not all use a four-stimulus judging set and, hence, do not all have a chance baseline of .25. The next most obvious candidate, the difference in each study between the hit rate observed and the hit rate expected under the null hypothesis, is also intuitively descriptive but is not appropriate for statistical analysis because not all differences between proportions that are equal are equally detectable (e.g., the power to detect the difference between .55 and .25 is different from the power to detect the difference between .50 and .20).

To provide a scale of equal detectability, Cohen (1988) devised the effect size index  $h$ , which involves an arcsine transformation on the proportions before calculation of their difference. Cohen's  $h$  is quite general and can assess the difference between any two proportions drawn from independent samples or between a single proportion and any specified hypothetical value. For the 28 studies examined in the meta-analyses,  $h$  was .28, with a 95% confidence interval from .11 to .45.

But because values of  $h$  do not provide an intuitively descriptive scale, Rosenthal and Rubin (1989; Rosenthal, 1991) have recently suggested a new index,  $\pi$ , which applies specifically to one-sample, multiple-choice data of the kind obtained in ganzfeld experiments. In particular,  $\pi$  expresses all hit rates as the proportion of hits that would have been obtained if there had been only two equally likely alternatives—essentially a coin flip. Thus,  $\pi$  ranges from 0 to 1, with .5 expected under the null hypothesis. The formula is

$$\pi = \frac{Pk - 1}{Pk - 2 + 1},$$

where  $P$  is the raw proportion of hits and  $k$  is the number of alternative choices available. Because  $\pi$  has such a straightforward intuitive interpretation, we use it (or its conversion back to an equivalent four-alternative hit rate) throughout this article whenever it is applicable.

For the 28 studies examined in the meta-analyses, the mean value of  $\pi$  was .62, with a 95% confidence interval from .55 to .69. This corresponds to a four-alternative hit rate of 35%, with a 95% confidence interval from 28% to 43%.

Cohen (1988, 1992) has also categorized effect sizes into small, medium, and large, with medium denoting an effect size that should be apparent to the naked eye of a careful observer. For a statistic such as  $\pi$ , which indexes the deviation of a pro-

portion from .5, Cohen considers .65 to be a medium effect size. A statistically unaided observer should be able to detect the bias of a coin that comes up heads on 65% of the trials. Thus, at .62, the psi ganzfeld effect size falls just short of Cohen's naked-eye criterion. From the phenomenology of the ganzfeld experimenter, the corresponding hit rate of 35% implies that he or she will see a subject obtain a hit approximately every third session rather than every fourth.

It is also instructive to compare the psi ganzfeld effect with the results of a recent medical study that sought to determine whether aspirin can prevent heart attacks (Steering Committee of the Physicians' Health Study Research Group, 1988). The study was discontinued after 6 years because it was already clear that the aspirin treatment was effective ( $p < .00001$ ) and it was considered unethical to keep the control group on placebo medication. The study was widely publicized as a major medical breakthrough. But despite its undisputed reality and practical importance, the size of the aspirin effect is quite small: Taking aspirin reduces the probability of suffering a heart attack by only .008. The corresponding effect size ( $h$ ) is .068, about one third to one fourth the size of the psi ganzfeld effect (Atkinson et al., 1993, p. 236; Uts, 1991b).

In sum, we believe that the psi ganzfeld effect is large enough to be of both theoretical interest and potential practical importance.

### Experimental Correlates of the Psi Ganzfeld Effect

We showed earlier that the technique of correlating variables with effect sizes across studies can help to assess whether methodological flaws might have produced artifactual positive outcomes. The same technique can be used more affirmatively to explore whether an effect varies systematically with conceptually relevant variations in experimental procedure. The discovery of such correlates can help to establish an effect as genuine, suggest ways of increasing replication rates and effect sizes, and enhance the chances of moving beyond the simple demonstration of an effect to its explanation. This strategy is only heuristic, however. Any correlates discovered must be considered quite tentative, both because they emerge from post hoc exploration and because they necessarily involve comparisons across heterogeneous studies that differ simultaneously on many interrelated variables, known and unknown. Two such correlates emerged from the meta-analyses of the psi ganzfeld effect.

*Single-versus multiple-image targets.* Although most of the 28 studies in the meta-analysis used single pictures as targets, 9 (conducted by three different investigators) used View Master stereoscopic slide reels that presented multiple images focused on a central theme. Studies using the View Master reels produced significantly higher hit rates than did studies using the single-image targets (50% vs. 34%), ( $T(26) = 2.22, p = .035$ , two-tailed).

*Sender-receiver pairing.* In 17 of the 28 studies, participants were free to bring in friends to serve as senders. In 8 studies, only laboratory-assigned senders were used. (Three studies used no sender.) Unfortunately, there is no record of how many participants in the former studies actually brought in friends. Nevertheless, those 17 studies (conducted by six different investigators) had significantly higher hit rates than did the studies

size: that used only laboratory-assigned senders (44% vs. 26%),  $t(23) = 2.39, p = .025$ , two-tailed.

#### *The Joint Communiqué*

After their published exchange in 1985, Hyman and Honorton agreed to contribute a joint communiqué to the subsequent discussion that was published in 1986. First, they set forth their areas of agreement and disagreement:

We agree that there is an overall significant effect in this data base that cannot reasonably be explained by selective reporting or multiple analysis. We continue to differ over the degree to which the effect constitutes evidence for psi, but we agree that the final verdict awaits the outcome of future experiments conducted by a broader range of investigators and according to more stringent standards. (Hyman & Honorton, 1986, p. 351)

They then spelled out in detail the "more stringent standards" they believed should govern future experiments. These standards included strict security precautions against sensory leakage, testing and documentation of randomization methods for selecting targets and sequencing the judging pool, statistical correction for multiple analyses, advance specification of the status of the experiment (e.g., pilot study or confirmatory experiment), and full documentation in the published report of the experimental procedures and the status of statistical tests (e.g., planned or post hoc).

#### *The National Research Council Report*

In 1988, the National Research Council (NRC) of the National Academy of Sciences released a widely publicized report commissioned by the U.S. Army that assessed several controversial technologies for enhancing human performance, including accelerated learning, neurolinguistic programming, mental practice, biofeedback, and parapsychology (Druckman & Swets, 1988; summarized in Swets & Bjork, 1990). The report's conclusion concerning parapsychology was quite negative: "The Committee finds no scientific justification from research conducted over a period of 130 years for the existence of parapsychological phenomena" (Druckman & Swets, 1988, p. 22).

An extended refutation strongly protesting the committee's treatment of parapsychology has been published elsewhere (Palmer et al., 1989). The pertinent point here is simply that the NRC's evaluation of the ganzfeld studies does not reflect an additional, independent examination of the ganzfeld database but is based on the same meta-analysis conducted by Hyman that we have discussed in this article.

Hyman chaired the NRC's Subcommittee on Parapsychology, and, although he had concurred with Honorton 2 years earlier in their joint communiqué that "there is an overall significant effect in this data base that cannot reasonably be explained by selective reporting or multiple analysis" (p. 351) and that "significant outcomes have been produced by a number of different investigators" (p. 352), neither of these points is acknowledged in the committee's report.

The NRC also solicited a background report from Harris and Rosenthal (1988a), which provided the committee with a comparative methodological analysis of the five controversial areas just listed. Harris and Rosenthal noted that, of these areas,

"only the Ganzfeld ESP studies [the only psi studies they evaluated] regularly meet the basic requirements of sound experimental design" (p. 53), and they concluded that

it would be implausible to entertain the null given the combined  $p$  from these 28 studies. Given the various problems or flaws pointed out by Hyman and Honorton . . . we might estimate the obtained accuracy rate to be about 1/3 . . . when the accuracy rate expected under the null is 1/4. (p. 51)<sup>3</sup>

#### The Autoganzfeld Studies

In 1983, Honorton and his colleagues initiated a new series of ganzfeld studies designed to avoid the methodological problems he and others had identified in earlier studies (Honorton, 1979; Kennedy, 1979). These studies complied with all of the detailed guidelines that he and Hyman were to publish later in their joint communiqué. The program continued until September 1989, when a loss of funding forced the laboratory to close. The major innovations of the new studies were computer control of the experimental protocol—hence the name *autoganzfeld*—and the introduction of videotaped film clips as target stimuli.

#### *Method*

The basic design of the autoganzfeld studies was the same as that described earlier<sup>4</sup>. A receiver and sender were sequestered in separate, acoustically isolated chambers. After a 14-min period of progressive relaxation, the receiver underwent ganzfeld stimulation while describing his or her thoughts and images aloud for 30 min. Meanwhile, the sender concentrated on a randomly selected target. At the end of the ganzfeld period, the receiver was shown four stimuli and, without knowing which of the four had been the target, rated each stimulus for its similarity to his or her mentation during the ganzfeld.

The targets consisted of 80 still pictures (static targets) and 80 short video segments complete with soundtracks (dynamic targets), all recorded on videocassette. The static targets included art prints, photographs, and magazine advertisements; the dynamic targets included excerpts of approximately 1-min duration from motion pictures, TV shows, and cartoons. The 160 targets were arranged in judging sets of four static or four dynamic targets each, constructed to minimize similarities among targets within a set.

*Target selection and presentation.* The VCR containing the taped targets was interfaced to the controlling computer, which selected the target and controlled its repeated presentation to the sender during the ganzfeld period, thus eliminating the need for a second experimenter to accompany the sender. After the ganzfeld period, the computer randomly sequenced the four-clip judging set and presented it to the receiver on a TV monitor for judging. The receiver used a computer game paddle to make his or her ratings on a 40-point scale that appeared on

<sup>3</sup> In a troubling development, the chair of the NRC Committee phoned Rosenthal and asked him to delete the parapsychology section of the paper (R. Rosenthal, personal communication, September 15, 1992). Although Rosenthal refused to do so, that section of the Harris-Rosenthal paper is nowhere cited in the NRC report.

<sup>4</sup> Because Honorton and his colleagues have complied with the Hyman-Honorton specification that experimental reports be sufficiently complete to permit others to reconstruct the investigator's procedures, readers who wish to know more detail than we provide here are likely to find whatever they need in the archival publication of these studies in the *Journal of Parapsychology* (Honorton et al., 1990).

the TV monitor after each clip was shown. The receiver was permitted to see each clip and to change the ratings repeatedly until he or she was satisfied. The computer then wrote these and other data from the session into a file on a floppy disk. At that point, the sender moved to the receiver's chamber and revealed the identity of the target to both the receiver and the experimenter. Note that the experimenter did not even know the identity of the four-clip judging set until it was displayed to the receiver for judging.

**Randomization.** The random selection of the target and sequencing of the judging set were controlled by a noise-based random number generator interfaced to the computer. Extensive testing confirmed that the generator was providing a uniform distribution of values throughout the full target range (1-160). Tests on the actual frequencies observed during the experiments confirmed that targets were, on average, selected uniformly from among the 4 clips within each judging set and that the 4 judging sequences used were uniformly distributed across sessions.

**Additional control features.** The receiver's and sender's rooms were sound-isolated, electrically shielded chambers with single-door access that could be continuously monitored by the experimenter. There was two-way intercom communication between the experimenter and the receiver but only one-way communication into the sender's room; thus, neither the experimenter nor the receiver could monitor events inside the sender's room. The archival record for each session includes an audiotape containing the receiver's mentation during the ganzfeld period and all verbal exchanges between the experimenter and the receiver throughout the experiment.

The automated ganzfeld protocol has been examined by several dozen parapsychologists and behavioral researchers from other fields, including well-known critics of parapsychology. Many have participated as subjects or observers. All have expressed satisfaction with the handling of security issues and controls.

Parapsychologists have often been urged to employ magicians as consultants to ensure that the experimental protocols are not vulnerable either to inadvertent sensory leakage or to deliberate cheating. Two "mentalists," magicians who specialize in the simulation of psi, have examined the autoganzfeld system and protocol. Ford Kross, a professional mentalist and officer of the mentalist's professional organization, the Psychic Entertainers Association, provided the following written statement: "In my professional capacity as a mentalist, I have reviewed Psychophysical Research Laboratories' automated ganzfeld system and found it to provide excellent security against deception by subjects" (personal communication, May, 1989).

Daryl J. Bem has also performed as a mentalist for many years and is a member of the Psychic Entertainers Association. As mentioned in the author note, this article had its origins in a 1983 visit he made to Honorton's laboratory, where he was asked to critically examine the research protocol from the perspective of a mentalist, a research psychologist, and a subject. Needless to say, this article would not exist if he did not concur with Ford Kross's assessment of the security procedures.

### Experimental Studies

Altogether, 100 men and 140 women participated as receivers in 354 sessions during the research program.<sup>5</sup> The participants ranged in age from 17 to 74 years ( $M = 37.3$ ,  $SD = 11.8$ ), with a mean formal education of 15.6 years ( $SD = 2.0$ ). Eight separate experimenters, including Honorton, conducted the studies.

The experimental program included three pilot and eight formal studies. Five of the formal studies used novice (first-time) participants who served as the receiver in one session each. The remaining three formal studies used experienced participants.

**Pilot studies.** Sample sizes were not present in the three pilot studies. Study 1 comprised 22 sessions and was conducted during the initial development and testing of the autoganzfeld system. Study 2 comprised

9 sessions testing a procedure in which the experimenter, rather than the receiver, served as the judge at the end of the session. Study 3 comprised 35 sessions and served as practice for participants who had completed the allotted number of sessions in the ongoing formal studies but who wanted additional ganzfeld experience. This study also included several demonstration sessions when TV film crews were present.

**Novice studies.** Studies 101-104 were each designed to test 50 participants who had had no prior ganzfeld experience; each participant served as the receiver in a single ganzfeld session. Study 104 included 16 of 20 students recruited from the Julliard School in New York City to test an artistically gifted sample. Study 105 was initiated to accommodate the overflow of participants who had been recruited for Study 104, including the 4 remaining Julliard students. The sample size for this study was set to 25, but only 6 sessions had been completed when the laboratory closed. For purposes of exposition, we divided the 56 sessions from Studies 104 and 105 into two parts: Study 104a/105(a) comprises the 36 non-Julliard participants, and Study 104b/105(b) comprises the 20 Julliard students.

**Study 201.** This study was designed to retest the most promising participants from the previous studies. The number of trials was set to 20, but only 7 sessions with 3 participants had been completed when the laboratory closed.

**Study 301.** This study was designed to compare static and dynamic targets. The sample size was set to 50 sessions. Twenty-five experienced participants each served as the receiver in 2 sessions. Unknown to the participants, the computer control program was modified to ensure that they would each have 1 session with a static target and 1 session with a dynamic target.

**Study 302.** This study was designed to examine a dynamic target set that had yielded a particularly high hit rate in the previous studies. The study involved experienced participants who had had no prior experience with this particular target set and who were unaware that only one target set was being sampled. Each served as the receiver in a single session. The design called for the study to continue until 15 sessions were completed with each of the targets, but only 25 sessions had been completed when the laboratory closed.

The 11 studies just described comprise all sessions conducted during the 6.5 years of the program. There is no "file drawer" of unreported sessions.

### Results

**Overall hit rate.** As in the earlier meta-analysis, receivers' ratings were analyzed by tallying the proportion of hits achieved and calculating the exact binomial probability for the observed number of hits compared with the chance expectation of .25. As noted earlier, 240 participants contributed 354 sessions. For reasons discussed later, Study 302 is analyzed separately, reducing the number of sessions in the primary analysis to 329.

As Table 1 shows, there were 106 hits in the 329 sessions, a hit rate of 32% ( $z = 2.89$ ,  $p = .002$ , one-tailed), with a 95% confidence interval from 30% to 35%. This corresponds to an effect size ( $\pi$ ) of .59, with a 95% confidence interval from .53 to .64.

Table 1 also shows that when Studies 104 and 105 are combined and re-divided into Studies 104/105(a) and 104/105(b), 9

<sup>5</sup> A recent review of the original computer files uncovered a duplicate record in the autoganzfeld database. This has now been eliminated, reducing by one the number of subjects and sessions. As a result, some of the numbers presented in this article differ slightly from those in Honorton et al. (1990).

Table 1  
Outcome by Study

Study	Study/subject description	N subjects	N trials	N hits	% hits	Effect size $\pi$	$z$
1	Pilot	19	22	8	36	.62	0.99
2	Pilot	4	9	3	33	.60	0.25
3	Pilot	24	35	10	29	.55	0.32
101	Novice	50	50	12	24	.47	-0.30
102	Novice	50	50	18	36	.63	
103	Novice	50	50	15	30	.55	0.67
104/105(a)	Novice	36	36	12	33	.60	0.97
104/105(b)	Juilliard sample	20	20	10	50	.75	2.20
201	Experienced	3	7	3	43	.69	0.69
301	Experienced	25	50	15	30	.66	0.67
302	Experienced	25	25	16	54*	.78*	3.04*
Overall (Studies 1-301)		240	329	106	32	.59	2.89

Note. All  $z$  scores are based on the exact binomial probability, with  $p = .25$  and  $q = .75$ .

\* Adjusted for response bias; the hit rate actually observed was 64%.

of the 10 studies yield positive effect sizes, with a mean effect size ( $\pi$ ) of .61,  $t(9) = 4.44$ ,  $p = .0008$ , one-tailed. This effect size is equivalent to a four-alternative hit rate of 34%. Alternatively, if Studies 104 and 105 are retained as separate studies, 9 of the 10 studies again yield positive effect sizes, with a mean effect size ( $\pi$ ) of .62,  $t(9) = 3.73$ ,  $p = .002$ , one-tailed. This effect size is equivalent to a four-alternative hit rate of 35% and is identical to that found across the 28 studies of the earlier meta-analysis.<sup>6</sup> Considered together, sessions with novice participants (Studies 101-105) yielded a statistically significant hit rate of 32.5% ( $p = .009$ ), which is not significantly different from the 31.6% hit rate achieved by experienced participants in Studies 201 and 301. And, finally, each of the eight experimenters also achieved a positive effect size, with a mean  $\pi$  of .60,  $t(7) = 3.44$ ,  $p = .005$ , one-tailed.

*The Juilliard sample.* There are several reports in the literature of a relationship between creativity or artistic ability and psi performance (Schmeidler, 1988). To explore this possibility in the ganzfeld setting, 10 male and 10 female undergraduates were recruited from the Juilliard School. Of these, 8 were music students, 10 were drama students, and 2 were dance students. Each served as the receiver in a single session in Study 104 or 105. As shown in Table 1, these students achieved a hit rate of 50% ( $p = .014$ ), one of the five highest hit rates ever reported for a single sample in a ganzfeld study. The musicians were particularly successful: 6 of the 8 (75%) successfully identified their targets ( $p = .004$ ; further details about this sample and their ganzfeld performance were reported in Schilitz & Honorton, 1992).

*Study size and effect size.* There is a significant negative correlation across the 10 studies listed in Table 1 between the number of sessions included in a study and the study's effect size ( $\pi$ ),  $r = -.64$ ,  $t(8) = 2.36$ ,  $p < .05$ , two-tailed. This is reminiscent of Hyman's discovery that the smaller studies in the original ganzfeld database were disproportionately likely to report statistically significant results. He interpreted this finding as evidence for a bias against the reporting of small studies that fail to achieve significant results. A similar interpretation cannot be

applied to the autoganzfeld studies, however, because there are no unreported sessions.

One reviewer of this article suggested that the negative correlation might reflect a decline effect in which earlier sessions of a study are more successful than later sessions. If there were such an effect, then studies with fewer sessions would show larger effect sizes because they would end before the decline could set in. To check this possibility, we computed point-biserial correlations between hits (1) or misses (0) and the session number within each of the 10 studies. All of the correlations hovered around zero; six were positive, four were negative, and the overall mean was .01.

An inspection of Table 1 reveals that the negative correlation derives primarily from the two studies with the largest effect sizes: the 20 sessions with the Juilliard students and the 7 sessions of Study 201, the study specifically designed to retest the most promising participants from the previous studies. Accordingly, it seems likely that the larger effect sizes of these two studies—and hence the significant negative correlation between the number of sessions and the effect size—reflect genuine performance differences between these two small, highly selected samples and other autoganzfeld participants.

*Study 302.* All of the studies except Study 302 randomly sampled from a pool of 160 static and dynamic targets. Study 302 sampled from a single, dynamic target set that had yielded a particularly high hit rate in the previous studies. The four film clips in this set consisted of a scene of a tidal wave from the movie *Clash of the Titans*, a high-speed sex scene from *A Clockwork Orange*, a scene of crawling snakes from a TV documentary, and a scene from a Bugs Bunny cartoon.

<sup>6</sup> As noted above, the laboratory was forced to close before three of the formal studies could be completed. If we assume that the remaining trials in Studies 105 and 201 would have yielded only chance results, this would reduce the overall  $z$  for the first 10 autoganzfeld studies from 2.89 to 2.76 ( $p = .003$ ). Thus, inclusion of the two incomplete studies does not pose an optional stopping problem. The third incomplete study, Study 302, is discussed below.

The experimental design called for this study to continue until each of the clips had served as the target 15 times. Unfortunately, the premature termination of this study at 25 sessions left an imbalance in the frequency with which each clip had served as the target. This means that the high hit rate observed (64%) could well be inflated by response biases.

As an illustration, water imagery is frequently reported by receivers in ganzfeld sessions, whereas sexual imagery is rarely reported. (Some participants probably are reluctant both to report sexual imagery and to give the highest rating to the sex-related clip.) If a video clip containing popular imagery (such as water) happens to appear as a target more frequently than a clip containing unpopular imagery (such as sex), a high hit rate might simply reflect the coincidence of those frequencies of occurrence with participants' response biases. And, as the second column of Table 2 reveals, the tidal wave clip did in fact appear more frequently as the target than did the sex clip. More generally, the second and third columns of Table 2 show that the frequency with which each film clip was ranked first closely matches the frequency with which each appeared as the target.

One can adjust for this problem by using the observed frequencies in these two columns to compute the hit rate expected if there were no psi effect. In particular, one can multiply each proportion in the second column by the corresponding proportion in the third column—yielding the joint probability that the clip was the target and that it was ranked first—and then sum across the four clips. As shown in the fourth column of Table 2, this computation yields an overall expected hit rate of 34.08%. When the observed hit rate of 64% is compared with this baseline, the effect size ( $h$ ) is .61. As shown in Table 1, this is equivalent to a four-alternative hit rate of 54%, or a  $\pi$  value of .78, and is statistically significant ( $z = 3.04, p = .0012$ ).

The psi effect can be seen even more clearly in the remaining columns of Table 2, which control for the differential popularity of the imagery in the clips by displaying how frequently each was ranked first when it was the target and how frequently it was ranked first when it was one of the control clips (decoys). As can be seen, each of the four clips was selected as the target relatively more frequently when it was the target than when it was a decoy, a difference that is significant for three of the four clips. On average, a clip was identified as the target 58% of the time when it was the target and only 14% of the time when it was a decoy.

*Dynamic versus static targets.* The success of Study 302 raises the question of whether dynamic targets are, in general, more effective than static targets. This possibility was also suggested by the earlier meta-analysis, which revealed that studies using multiple-image targets (View Master stereoscopic slide reels) obtained significantly higher hit rates than did studies using single-image targets. By adding motion and sound, the video clips might be thought of as high-tech versions of the View Master reels.

The 10 autoganzfeld studies that randomly sampled from both dynamic and static target pools yielded 164 sessions with dynamic targets and 165 sessions with static targets. As predicted, sessions using dynamic targets yielded significantly more hits than did sessions using static targets (37% vs. 27%; Fisher's exact  $p < .04$ ).

*Sender-receiver pairing.* The earlier meta-analysis revealed that studies in which participants were free to bring in friends

Table 2  
Study 302: Expected Hit Rate and Proportion of Sessions in Which Each Video Clip Was Ranked First When It Was a Target and When It Was a Decoy

Video clip	Relative frequency as target	Relative frequency of first place ranking	Expected hit rate (%)	Ranked first when target	Ranked first when decoy	Difference	Fisher's exact $p$
Tidal wave	.28 (7/25)	.24 (6/25)	6.72	.57 (4/7)	.11 (2/18)	.46	.032
Snakes	.12 (3/25)	.12 (3/25)	1.44	.67 (2/3)	.05 (1/22)	.62	.029
Sex scene	.16 (4/25)	.08 (2/25)	1.28	.25 (1/4)	.05 (1/21)	.20	.300
Bugs Bunny	.44 (11/25)	.56 (14/25)	24.64	.82 (9/11)	.36 (5/14)	.46	.027
Overall			34.08	.58	.14	.44	

to serve as senders produced significantly higher hit rates than studies that used only laboratory-assigned senders. As noted, however, there is no record of how many of the participants in the former studies actually did bring in friends. Whatever the case, sender-receiver pairing was not a significant correlate of psi performance in the autoganzfeld studies: The 197 sessions in which the sender and receiver were friends did not yield a significantly higher proportion of hits than did the 132 sessions in which they were not (35% vs. 29%; Fisher's exact  $p = .28$ ).

*Correlations between receiver characteristics and psi performance.* Most of the autoganzfeld participants were strong believers in psi: On a 7-point scale ranging from *strong disbelief in psi* (1) to *strong belief in psi* (7), the mean was 6.2 ( $SD = 1.03$ ); only 2 participants rated their belief in psi below the midpoint of the scale. In addition, 88% of the participants reported personal experiences suggestive of psi, and 80% had some training in meditation or other techniques involving internal focus of attention.

All of these appear to be important variables. The correlation between belief in psi and psi performance is one of the most consistent findings in the parapsychological literature (Palmer, 1978). And, within the autoganzfeld studies, successful performance of novice (first-time) participants was significantly predicted by reported personal psi experiences, involvement with meditation or other mental disciplines, and high scores on the Feeling and Perception factors of the Myers-Briggs Type Inventory (Honorton, 1992; Honorton & Schechter, 1987; Myers & McCaulley, 1985). This recipe for success has now been independently replicated in another laboratory (Broughton, Kanthamam, & Khilji, 1990).

The personality trait of extraversion is also associated with better psi performance. A meta-analysis of 60 independent studies with nearly 3,000 subjects revealed a small but reliable positive correlation between extraversion and psi performance, especially in studies that used free-response methods of the kind used in the ganzfeld experiments (Honorton, Ferrari, & Bern, 1992). Across 14 free-response studies conducted by four independent investigators, the correlation for 612 subjects was  $.20$  ( $z = 4.82$ ,  $p = 1.5 \times 10^{-6}$ ). This correlation was replicated in the autoganzfeld studies, in which extraversion scores were available for 218 of the 240 subjects,  $r = .18$ ,  $t(216) = 2.67$ ,  $p = .004$ , one-tailed.

Finally, there is the strong psi performance of the Julliard students, discussed earlier, which is consistent with other studies in the parapsychological literature suggesting a relationship between successful psi performance and creativity or artistic ability.

## Discussion

Earlier in this article, we quoted from the abstract of the Hyman-Honorton (1986) communiqué: "We agree that the final verdict awaits the outcome of future experiments conducted by a broader range of investigators and according to more stringent standards" (p. 351). We believe that the "stringent standards" requirement has been met by the autoganzfeld studies. The results are statistically significant and consistent with those in the earlier database. The mean effect size is quite respectable in comparison with other controversial research areas of human

performance (Harris & Rosenthal, 1988a). And there are reliable relationships between successful psi performance and conceptually relevant experimental and subject variables, relationships that also replicate previous findings. Hyman (1991) has also commented on the autoganzfeld studies: "Honorton's experiments have produced intriguing results. If . . . independent laboratories can produce similar results with the same relationships and with the same attention to rigorous methodology, then parapsychology may indeed have finally captured its elusive quarry" (p. 392).

## Issues of Replication

As Hyman's comment implies, the autoganzfeld studies by themselves cannot satisfy the requirement that replications be conducted by a "broader range of investigators." Accordingly, we hope the findings reported here will be sufficiently provocative to prompt others to try replicating the psi ganzfeld effect.

We believe that it is essential, however, that future studies comply with the methodological, statistical, and reporting standards set forth in the joint communiqué and achieved by the autoganzfeld studies. It is not necessary for studies to be as automated or as heavily instrumented as the autoganzfeld studies to satisfy the methodological guidelines, but they are still likely to be labor intensive and potentially expensive.<sup>7</sup>

## Statistical Power and Replication

Would-be replicators also need to be reminded of the power requirements for replicating small effects. Although many academic psychologists do not believe in psi, many apparently do believe in miracles when it comes to replication. Tversky and Kahneman (1971) posed the following problem to their colleagues at meetings of the Mathematical Psychology Group and the American Psychological Association:

Suppose you have run an experiment on 20 subjects and have obtained a significant result which confirms your theory ( $z = 2.23$ ,  $p < .05$ , two-tailed). You now have cause to run an additional group of 10 subjects. What do you think the probability is that the results will be significant, by a one-tailed test, separately for this group? (p. 105)

The median estimate was .85, with 9 of 10 respondents providing an estimate greater than .60. The correct answer is approximately .48.

As Rosenthal (1990) has warned: "Given the levels of statistical power at which we normally operate, we have no right to expect the proportion of significant results that we typically do expect, even if in nature there is a very real and very important effect" (p. 16). In this regard, it is again instructive to consider the medical study that revealed a highly significant effect of aspirin on the incidence of heart attacks. The study monitored more than 22,000 subjects. Had the investigators monitored 3,000 subjects, they would have had less than an even chance of

<sup>7</sup> As the closing of the autoganzfeld laboratory exemplifies, it is also difficult to obtain funding for psi research. The traditional, peer-reviewed sources of funding familiar to psychologists have almost never funded proposals for psi research. The widespread skepticism of psychologists toward psi is almost certainly a contributing factor.

finding a conventionally significant effect. Such is life with small effect sizes.

Given its larger effect size, the prospects for successfully replicating the psi ganzfeld effect are not quite so daunting, but they are probably still grimmer than intuition would suggest. If the true hit rate is in fact about 34% when 25% is expected by chance, then an experiment with 30 trials (the mean for the 28 studies in the original meta-analysis) has only about 1 chance in 6 of finding an effect significant at the .05 level with a one-tailed test. A 50-trial experiment boosts that chance to about 1 in 3. One must escalate to 100 trials to come close to the break-even point, at which one has a 50-50 chance of finding a statistically significant effect (Ufrs, 1986). (Recall that only 2 of the 11 autoganzfeld studies yielded results that were individually significant at the conventional .05 level.) Those who require that a psi effect be statistically significant every time before they will seriously entertain the possibility that an effect really exists know not what they ask.

### Significance Versus Effect Size

The preceding discussion is unduly pessimistic, however, because it perpetuates the tradition of worshipping the significance level. Regular readers of this journal are likely to be familiar with recent arguments imploring behavioral scientists to overcome their slavish dependence on the significance level as the ultimate measure of virtue and instead to focus more of their attention on effect sizes: "Surely, God loves the .06 nearly as much as the .05" (Rosnow & Rosenthal, 1989, p. 1277). Accordingly, we suggest that achieving a respectable effect size with a methodologically tight ganzfeld study would be a perfectly welcome contribution to the replication effort, no matter how untenurable the  $p$  level renders the investigator.

Career consequences aside, this suggestion may seem quite counterintuitive. Again, Tversky and Kahneman (1971) have provided an elegant demonstration. They asked several of their colleagues to consider an investigator who runs 15 subjects and obtains a significant  $t$  value of 2.46. Another investigator attempts to duplicate the procedure with the same number of subjects and obtains a result in the same direction but with a nonsignificant value of  $t$ . Tversky and Kahneman then asked their colleagues to indicate the highest level of  $t$  in the replication study they would describe as a failure to replicate. The majority of their colleagues regarded  $t = 1.70$  as a failure to replicate. But if the data from two such studies ( $t = 2.46$  and  $t = 1.70$ ) were pooled, the  $t$  for the combined data would be about 10 (assuming equal variances):

Thus, we are faced with a paradoxical state of affairs, in which the same data that would increase our confidence in the finding when viewed as part of the original study, shake our confidence when viewed as an independent study. (Tversky & Kahneman, 1971, p. 108)

Such is the iron grip of the arbitrary .05. Pooling the data, of course, is what meta-analysis is all about. Accordingly, we suggest that two or more laboratories could collaborate in a ganzfeld replication effort by conducting independent studies and pooling them in meta-analytic fashion, what one might call a time meta-analysis. (Each investigator could then claim the  $ed p$  level for his or her own curriculum vitae.)

### Maximizing Effect Size

Rather than buying or borrowing larger sample sizes, those who seek to replicate the psi ganzfeld effect might find it more intellectually satisfying to attempt to maximize the effect size by attending to the variables associated with successful outcomes. Thus, researchers who wish to enhance the chances of successful replication should use dynamic rather than static targets. Similarly, we advise using participants with the characteristics we have reported to be correlated with successful performance. Random college sophomores enrolled in introductory psychology do not constitute the optimal subject pool. Finally, we urge ganzfeld researchers to read carefully the detailed description of the warm social ambience that Honorton et al. (1990) sought to create in the autoganzfeld laboratory. We believe that the social climate created in psi experiments is a critical determinant of their success or failure.

### The Problem of "Other" Variables

This caveat about the social climate of the ganzfeld experiment prompted one reviewer of this article to worry that this provided "an escape clause" that weakens the falsifiability of the psi hypothesis: "Until Bem and Honorton can provide operational criteria for creating a warm social ambience, the failure of an experiment with otherwise adequate power can always be dismissed as due to a lack of warmth."

Alas, it is true; we devoutly wish it were otherwise. But the operation of unknown variables in moderating the success of replications is a fact of life in all of the sciences. Consider, for example, an earlier article in this journal by Spence (1964). He reviewed studies testing the straightforward derivation from Hullian learning theory that high-anxiety subjects should condition more strongly than low-anxiety subjects. This hypothesis was confirmed 94% of the time in Spence's own laboratory at the University of Iowa but only 63% of the time in laboratories at other universities. In fact, Kimble and his associates at Duke University and the University of North Carolina obtained results in the opposite direction in two of three experiments.

In searching for a post hoc explanation, Spence (1964) noted that "a deliberate attempt was made in the Iowa studies to provide conditions in the laboratory that might elicit some degree of emotionality. Thus, the experimenter was instructed to be impersonal and quite formal . . . and did not try to put [subjects] at ease or allay any expressed fears" (pp. 135-136). Moreover, he pointed out, his subjects sat in a denial chair, whereas Kimble's subjects sat in a secretarial chair. Spence even considered "the possibility that cultural backgrounds of southern and northern students may lead to a difference in the manner in which they respond to the different items in the [Manifest Anxiety] scale" (p. 136). If this was the state of affairs in an area of research as well established as classical conditioning, then the suggestion that the social climate of the psi laboratory might affect the outcome of ganzfeld experiments in ways not yet completely understood should not be dismissed as a devious attempt to provide an escape clause in case of replication failure.

The best the original researchers can do is to communicate as completely a knowledge of the experimental conditions as possible in an attempt to anticipate some of the relevant moderating

variables. Ideal research, however, the dures provide rent knowledge "operational

Up to this empirical me should establish not, before a moment that transfer here.

### The Psychol

In attempti cally begun w derlying meet psychological that target inf lus that is enc formation-pr mances shoul in psycholog iced in the mox first place.

The ganzfe ganzfeld proc mediated infc is normally n "noise." Acco noise ratio sh ated informat a large and di states of conse and, of course Alternative states) may be cepting alien strains on the more divergen ritual that par ford, 1987). A one to choose; model remain: The larger,

that attempt it over static targ involve more ; internal schem are more emot specified ways. yond the simpl theory-based d have involved

variables. Ideally, this might include direct training by the original researchers or videotapes of actual sessions. Lacking these, however, the detailed description of the autoganzfeld procedures provided by Honorton et al. (1990) comes as close as current knowledge permits in providing for other researchers the "operational criteria for creating a warm social ambiance."

#### Theoretical Considerations

Up to this point, we have confined our discussion to strictly empirical matters. We are sympathetic to the view that one should establish the existence of a phenomenon, anomalous or not, before attempting to explain it. So let us suppose for the moment that we have a genuine anomaly of information transfer here. How can it be understood or explained?

#### *The Psychology of Psi*

In attempting to understand psi, parapsychologists have typically begun with the working assumption that, whatever its underlying mechanisms, it should behave like other, more familiar psychological phenomena. In particular, they typically assume that target information behaves like an external sensory stimulus that is encoded, processed, and experienced in familiar information-processing ways. Similarly, individual psi performances should covary with experimental and subject variables in psychologically sensible ways. These assumptions are embodied in the model of psi that motivated the ganzfeld studies in the first place.

*The ganzfeld procedure.* As noted in the introduction, the ganzfeld procedure was designed to test a model in which psi-mediated information is conceptualized as a weak signal that is normally masked by internal somatic and external sensory "noise." Accordingly, any technique that raises the signal-to-noise ratio should enhance a person's ability to detect psi-mediated information. This noise-reduction model of psi organizes a large and diverse body of experimental results, particularly those demonstrating the psi-conductive properties of altered states of consciousness such as meditation, hypnosis, dreaming, and, of course, the ganzfeld itself (Rao & Palmer, 1987).

Alternative theories propose that the ganzfeld (and altered states) may be psi conductive because it lowers resistance to accepting alien imagery, diminishes rational or contextual constraints on the encoding or reporting of information, stimulates more divergent thinking, or even just serves as a placebo-like ritual that participants perceive as being psi conductive (Stanford, 1987). At this point, there are no data that would permit one to choose among these alternatives, and the noise-reduction model remains the most widely accepted.

*The target.* There are also a number of plausible hypotheses that attempt to account for the superiority of dynamic targets over static targets. Dynamic targets contain more information, involve more sensory modalities, evoke more of the receiver's internal schemata, are more lifelike, have a narrative structure, are more emotionally evocative, and are "richer" in other, unspecified ways. Several psi researchers have attempted to go beyond the simple dynamic-static dichotomy to more refined or theory-based definitions of a good target. Although these efforts have involved examining both psychological and physical prop-

erties of targets, there is as yet not much progress to report (Delano, 1990).

*The receiver.* Some of the subject characteristics associated with good psi performance also appear to have psychologically straightforward explanations. For example, garden-variety motivational explanations seem sufficient to account for the relatively consistent finding that those who believe in psi perform significantly better than those who do not. (Less straightforward, however, would be an explanation for the frequent finding that nonbelievers actually perform significantly worse than chance [Broughton, 1991, p. 109].)

The superior psi performance of creative or artistically gifted individuals—such as the Juilliard students—may reflect individual differences that parallel some of the hypothesized effects of the ganzfeld mentioned earlier: Artistically gifted individuals may be more receptive to alien imagery, be better able to transcend rational or contextual constraints in the encoding or reporting of information, or be more divergent in their thinking. It has also been suggested that both artistic and psi abilities might be rooted in superior right-brain functioning.

The observed relationship between extraversion and psi performance has been of theoretical interest for many years. Eysenck (1966) reasoned that extraverts should perform well in psi tasks because they are easily bored and respond favorably to novel stimuli. In a setting such as the ganzfeld, extraverts may become "stimulus starved" and thus may be highly sensitive to any stimulation, including weak incoming psi information. In contrast, introverts would be more inclined to entertain themselves with their own thoughts and thus continue to mask psi information despite the diminished sensory input. Eysenck also speculated that psi might be a primitive form of perception antedating cortical developments in the course of evolution, and, hence, cortical arousal might suppress psi functioning. Because extraverts have a lower level of cortical arousal than introverts, they should perform better in psi tasks (the evolutionary biology of psi has also been discussed by Broughton, 1991, pp. 347-352).

But there are more mundane possibilities. Extraverts might perform better than introverts simply because they are more relaxed and comfortable in the social setting of the typical psi experiment (e.g., the "warm social ambiance" of the autoganzfeld studies). This interpretation is strengthened by the observation that introverts outperformed extraverts in a study in which subjects had no contact with an experimenter but worked alone at home with materials they received in the mail (Schmidt & Schlitz, 1989). To help decide among these interpretations, ganzfeld experimenters have begun to use the extraversion scale of the NEO Personality Inventory (Costa & McCrae, 1992), which assesses six different facets of the extraversion-introversion factor.

*The sender.* In contrast to this information about the receiver in psi experiments, virtually nothing is known about the characteristics of a good sender or about the effects of the sender's relationship with the receiver. As has been shown, the initial suggestion from the meta-analysis of the original ganzfeld database that psi performance might be enhanced when the sender and receiver are friends was not replicated at a statistically significant level in the autoganzfeld studies.

A number of parapsychologists have entertained the more



- ld cease  
t of our  
yes and
- Reflections on Bell's theorem.* Notre Dame, IN: University of Notre Dame Press.
- Dawson, R. (1991). Comment. *Statistical Science*, 6, 382-385.
- DeJano, D. L. (1990). Approaches to the target: A time for reevaluation. In L. A. Henkel & J. Palmer (Eds.), *Research in parapsychology* /1989 (pp. 89-92). Metuchen, NJ: Scarecrow Press.
- Dingwall, E. J. (Ed.). (1968). *Abnormal hypnotic phenomena* (4 vols.). London: Churchill.
- Druckman, D., & Swets, J. A. (Eds.). (1988). *Enhancing human performance: Issues, theories, and techniques*. Washington, DC: National Academy Press.
- Eysenck, H. J. (1966). Personality and extra-sensory perception. *Journal of the Society for Psychological Research*, 44, 55-71.
- Gilovich, T. (1991). *How we know what isn't so: The fallibility of human reason in everyday life*. New York: Free Press.
- Green, C. E. (1960). Analysis of spontaneous cases. *Proceedings of the Society for Psychological Research*, 53, 97-161.
- Harris, M. J., & Rosenthal, R. (1988a). *Human performance research: An overview*. Washington, DC: National Academy Press.
- Harris, M. J., & Rosenthal, R. (1988b). *Postscript to "Human performance research: An overview"*. Washington, DC: National Academy Press.
- Herbert, N. (1987). *Quantum reality: Beyond the new physics*. Garden City, NY: Anchor Books.
- Honorton, C. (1969). Relationship between EEG alpha activity and ESP card-guessing performance. *Journal of the American Society for Psychological Research*, 63, 365-374.
- Honorton, C. (1977). Psi and internal attention states. In B. B. Wolman (Ed.), *Handbook of parapsychology* (pp. 435-472). New York: Van Nostrand Reinhold.
- Honorton, C. (1979). Methodological issues in free-response experiments. *Journal of the American Society for Psychological Research*, 73, 381-394.
- Honorton, C. (1985). Meta-analysis of psi ganzfeld research: A response to Hyman. *Journal of Parapsychology*, 49, 51-91.
- Honorton, C. (1992, August). *The ganzfeld novice: Four predictors of initial ESP performance*. Paper presented at the 35th annual convention of the Parapsychological Association, Las Vegas, NV.
- Honorton, C., Berger, R. E., Varvoglis, M. P., Quant, M., Derr, P., Schechter, E. I., & Ferrari, D. C. (1990). Psi communication in the ganzfeld: Experiments with an automated testing system and a comparison with a meta-analysis of earlier studies. *Journal of Parapsychology*, 54, 99-139.
- Honorton, C., Ferrari, D. C., & Bem, D. J. (1992). Extraversion and ESP performance: Meta-analysis and a new confirmation. In L. A. Henkel & G. R. Schneider (Eds.), *Research in parapsychology* 1990 (pp. 35-38). Metuchen, NJ: Scarecrow Press.
- Honorton, C., & Harper, S. (1974). Psi-mediated imagery and ideation in an experimental procedure for regulating perceptual input. *Journal of the American Society for Psychological Research*, 68, 156-168.
- Honorton, C., & Schechter, E. I. (1987). Ganzfeld target retrieval with an automated testing system: A model for initial ganzfeld success. In D. B. Weiner & R. D. Nelson (Eds.), *Research in parapsychology* /1986 (pp. 36-39). Metuchen, NJ: Scarecrow Press.
- Hyman, R. (1985). The ganzfeld psi experiment: A critical appraisal. *Journal of Parapsychology*, 49, 3-49.
- Hyman, R. (1991). Comment. *Statistical Science*, 6, 389-392.
- Hyman, R., & Honorton, C. (1986). A joint communiqué: The psi ganzfeld controversy. *Journal of Parapsychology*, 50, 351-364.
- Kennedy, J. E. (1979). Methodological problems in free-response ESP experiments. *Journal of the American Society for Psychological Research*, 73, 1-15.
- Mezger, W. (1930). Optische Untersuchungen am Ganzfeld: II. Zur Phänomenologie des homogenen Ganzfelds [Optical investigation of the Ganzfeld: II Toward the phenomenology of the homogeneous Ganzfeld]. *Psychologische Forschung*, 13, 6-29.
- Morris, R. L. (1991). Comment. *Statistical Science*, 6, 393-395.
- Myers, I. B., & McCaulley, M. H. (1985). *Manual: A guide to the development and use of the Myers-Briggs Type Indicator*. Palo Alto, CA: Consulting Psychologists Press.
- Nisbett, R. E., & Ross, L. (1980). *Human inference: Strategies and shortcomings of social judgment*. Englewood Cliffs, NJ: Prentice Hall.
- Palmer, J. (1978). Extrasensory perception: Research findings. In S. Krippner (Ed.), *Advances in parapsychological research* (Vol. 2, pp. 59-243). New York: Plenum Press.
- Palmer, J. A., Honorton, C., & Uts, J. (1989). Reply to the National Research Council Study on Parapsychology. *Journal of the American Society for Psychological Research*, 83, 31-49.
- Parker, A. (1975). Some findings relevant to the change in state hypothesis. In J. D. Morris, W. G. Roll, & R. L. Morris (Eds.), *Research in parapsychology* 1974 (pp. 40-42). Metuchen, NJ: Scarecrow Press.
- Parker, A. (1978). A holistic methodology in psi research. *Parapsychology Review*, 9, 1-6.
- Prasad, J., & Stevenson, I. (1968). A survey of spontaneous psychical experiences in school children of Uttar Pradesh, India. *International Journal of Parapsychology*, 10, 241-261.
- Rao, K. R., & Palmer, J. (1987). The anomaly called psi: Recent research and criticism. *Behavioral and Brain Sciences*, 10, 539-551.
- Rhine, J. B., & Pratt, J. G. (1954). A review of the Pearce-Pratt distance series of ESP tests. *Journal of Parapsychology*, 18, 165-177.
- Rhine, L. E. (1962). Psychological processes in ESP experiences. I. Waking experiences. *Journal of Parapsychology*, 26, 88-111.
- Roig, M., Icochea, H., & Guzzucoli, A. (1991). Coverage of parapsychology in introductory psychology textbooks. *Teaching of Psychology*, 18, 157-160.
- Rosenthal, R. (1978). Combining results of independent studies. *Psychological Bulletin*, 85, 185-193.
- Rosenthal, R. (1979). The "file drawer problem" and tolerance for null results. *Psychological Bulletin*, 86, 638-641.
- Rosenthal, R. (1990). Replication in behavioral research. *Journal of Social Behavior and Personality*, 5, 1-30.
- Rosenthal, R. (1991). *Meta-analytic procedures for social research* (Rev. ed.). Newbury Park, CA: Sage.
- Rosenthal, R., & Rubin, D. B. (1989). Effect size estimation for one-sample multiple-choice-type data: Design, analysis, and meta-analysis. *Psychological Bulletin*, 106, 332-337.
- Rosnow, R. L., & Rosenthal, R. (1989). Statistical procedures and the justification of knowledge in psychological science. *American Psychologist*, 44, 1276-1284.
- Sannwald, G. (1959). Statistische Untersuchungen an Spontanphänomene [Statistical investigation of spontaneous phenomena]. *Zeitschrift für Parapsychologie und Grenzgebiete der Psychologie*, 3, 59-71.
- Saunders, D. R. (1985). On Hyman's factor analyses. *Journal of Parapsychology*, 49, 86-88.
- Schechter, E. I. (1984). Hypnotic induction vs. control conditions: Illustrating an approach to the evaluation of replicability in parapsychology. *Journal of the American Society for Psychological Research*, 78, 1-27.
- Schultz, M. J., & Honorton, C. (1992). Ganzfeld psi performance within an artistically gifted population. *Journal of the American Society for Psychological Research*, 86, 83-98.
- Schmeidler, G. R. (1988). *Parapsychology and psychology: Matches and mismatches*. Jefferson, NC: McFarland.
- Schmidt, H., & Schultz, M. J. (1989). A large scale pilot PK experiment with prerecorded random events. In L. A. Henkel & R. E. Berger (Eds.), *Research in parapsychology* 1988 (pp. 6-10). Metuchen, NJ: Scarecrow Press.

- Spence, K. W. (1964). Anxiety (drive) level and performance in eyelid conditioning. *Psychological Bulletin*, 61, 129-139.
- Stanford, R. G. (1987). Ganzfeld and hypnotic-induction procedures in ESP research: Toward understanding their success. In S. Krippner (Ed.), *Advances in parapsychological research* (Vol. 5, pp. 39-76). Jefferson, NC: McFarland.
- Steering Committee of the Physicians' Health Study Research Group. (1988). Preliminary report: Findings from the aspirin component of the ongoing Physicians' Health Study. *New England Journal of Medicine*, 318, 262-264.
- Sterling, T. C. (1959). Publication decisions and their possible effects on inferences drawn from tests of significance—or vice versa. *Journal of the American Statistical Association*, 54, 30-34.
- Stokes, D. M. (1987). Theoretical parapsychology. In S. Krippner (Ed.), *Advances in parapsychological research* (Vol. 5, pp. 77-189). Jefferson, NC: McFarland.
- Swets, J. A., & Bjork, R. A. (1990). Enhancing human performance: An evaluation of "new age" techniques considered by the U.S. Army. *Psychological Science*, 1, 85-96.
- Tversky, A., & Kahneman, D. (1971). Belief in the law of small numbers. *Psychological Bulletin*, 2, 105-110.
- Ullman, M., Krippner, S., & Vaughan, A. (1973). *Dream telepathy*. New York: Macmillan.
- Uts, J. (1986). The ganzfeld debate: A statistician's perspective. *Journal of Parapsychology*, 50, 393-402.
- Uts, J. (1991a). Rejoinder. *Statistical Science*, 6, 396-403.
- Uts, J. (1991b). Replication and meta-analysis in parapsychology. *Statistical Science*, 6, 363-378.
- Wagner, M. W., & Monnet, M. (1979). Attitudes of college professors toward extrasensory perception. *Zetetic Scholar*, 5, 7-17.

Received September 28, 1992

Revision received March 10, 1993

Accepted March 14, 1993 ■

### Call for Nominations

The Publications and Communications Board has opened nominations for the editorships of *Behavioral Neuroscience*, the *Journal of Experimental Psychology: General*, and the *Journal of Experimental Psychology: Learning, Memory, and Cognition* for the years 1996-2001. Larry R. Squire, PhD, Earl Hunt, PhD, and Keith Rayner, PhD, respectively, are the incumbent editors. Candidates must be members of APA and should be available to start receiving manuscripts in early 1995 to prepare for issues published in 1996. Please note that the P&C Board encourages participation by members of underrepresented groups in the publication process and would particularly welcome such nominees. To nominate candidates, prepare a statement of one page or less in support of each candidate.

- For *Behavioral Neuroscience*, submit nominations to J. Bruce Overmier, PhD, Elliott Hall—Psychology, University of Minnesota, 75 East River Road, Minneapolis, MN 55455 or to psyjbo@vx.cis.umn.edu. Other members of the search committee are Norman Adler, PhD, Evelyn Sainoff, PhD, and Richard F. Thompson, PhD.
- For the *Journal of Experimental Psychology: General*, submit nominations to Howard E. Egeth, PhD, Chair, JEP: General Search, Department of Psychology, Johns Hopkins University, Charles & 34th Streets, Baltimore, MD 21218, to egeth@jhuvm.binet, or to fax number 410-516-4478. Other members of the search committee are Donald S. Blough, PhD, Martha Farah, PhD, and Edward E. Smith, PhD.
- For the *Journal of Experimental Psychology: Learning, Memory, and Cognition*, submit nominations to Donna M. Gelfand, PhD, Dean, Social and Behavioral Science, 205 Osh, University of Utah, Salt Lake City, UT 84112-1102 or to fax number 801-585-5081. Other members of the search committee are Marcia Johnson, PhD, Michael Posner, PhD, Henry L. Roediger III, PhD, and Richard M. Shiffrin, PhD.

First review of nominations will begin December 15, 1993.