

Tanulmány

Csilla Rákosi

On the evaluation of psycholinguistic experiments on metaphor

Part I: The metatheoretical background

Abstract

Psycholinguistic research into metaphor is characterised by contradictory and often controversial experimental results. Thus, the question is, how to decide when an experiment yields plausible experimental data and when it is unreliable as a data source. On the basis of a model of psycholinguistic experiments, it is proposed that experiments should be viewed as cyclic open processes. This means that the plausibility of the statements related to different stages of the experimental process is revised again and again during the elaboration and conduct of the experiment, as well as during its evaluation. Accordingly, the analysis and evaluation of experiments is nothing else than the continuation of the experimental process by new plausible argumentation cycles, and, if possible, the elaboration of proposals for its resumption by new experimental cycles.

Keywords: psycholinguistic experiments, experiments on metaphor processing, philosophy of science, evaluation of experiments

1 Introduction

As is well-known, psycholinguistic research into metaphor is characterised by a considerable diversity of approaches and experimental methods as well as contradictory and often controversial experimental results. Raymond W. Gibbs puts forward a two-step diagnosis of this situation. First, he claims that “psycholinguistic experiments may be [...] inherently flawed as a scientific enterprise” (Gibbs 2013: 45). Second, he raises the hypothesis that with the help of his alternative metascientific model, it is possible to “push metaphor scholars closer to thinking and practices seen in more mature scientific disciplines” (Gibbs 2013: 52).

In contrast, in Dirk Geeraerts’s view, psycholinguistic experiments apply feasible, well-established procedures providing completely reliable experimental results:

[...] there is a common, commonly accepted way in psycholinguistics of settling theoretical disputes: experimentation. Given a number of conditions, experimental results decide between competing analyses, and psycholinguists predominantly accept the experimental paradigm as the cornerstone of their discipline. (Geeraerts 2006: 26)

Hasson and Giora (2007) take another route: they provide us with a comprehensive overview of the experimental methods applied in cognitive linguistics, summarising their rationale and identifying their possible weak points. Their list can be profitably complemented with Keenan et al.’s (1990), Haberlandt’s (1994) and Kaiser’s (2013) considerations. This combined inventory, however, still cannot be regarded as a system of methodological guidelines, mainly due

to the circumstance that all three papers focus on the detailed characterisation of the basic hypotheses and working mechanism of the different experimental methods. Therefore, they provide neither a systematic nor an exhaustive typology of errors but discuss the most typical problems related to the different types of experiments.

This disagreement might motivate a twofold strategy. Namely, metascientific reflection on the nature and limits of psycholinguistic experiments should be based on the continuous comprehension and adjustment of *insights gained by philosophers of science studying experiments in science* (i.e., a model of scientific experiments in general) on the one hand, and the *reflection on the research activities of linguists while working with experiments* (that is, criteria related to the experimental methods used in linguistics, in particular), on the other. Both components are vital. First, as Rákosi (2014) argues, linguists often confuse workable and generally applied norms of natural sciences with outmoded and untenable tenets of the standard view of the analytical philosophy of science.¹ Second, contemporary philosophy of science rejects the idea of providing general, uniform norms for scientific theorising and experimenting. Instead, research practice is studied carefully and closely, and methodological rules or norms are held to be *field-sensitive* and put into a *historical context*.

Experiments involve many potential sources of error and undetected possibilities. Therefore, it would be vital to take the fallibility of experiments seriously and search for means which enable us to reduce it. One of the most effective ways of achieving this could be *replications*. Accordingly, this paper will raise and test the following hypothesis:

- (H) (a) The structure of psycholinguistic experiments basically *corresponds to the structure of experiments in the natural sciences*, although the role and impact of certain stages diverge.
- (b) Psycholinguistic experiments should be viewed as *cyclic open processes*. This means that the plausibility of the statements related to different stages of the experimental process should be re-evaluated again and again during the elaboration and conduct of the experiment, as well as during its evaluation.
- (c) The analysis and evaluation of experiments is nothing else than *the continuation of the experimental process by new plausible argumentation cycles*, and, if possible, the elaboration of *proposals for its continuation by new experimental cycles*.
- (d) The conduct of the proposed revised version(s) of the original experiment and the comparison of the results obtained may lead to *more elaborated experiments and more reliable experimental data*.
- (e) Through the *replications, strict revisions and improvements*, experiments should become *collective works of a research field* and not private affairs of single minds.
- (f) This also means that it is not only the “discovery phase” (the experimental process) that contains justification but the “justification phase” (the evaluation of experiments) also involves discovery.

As for the structure of this paper, Section 2 will present a metascientific model of psycholinguistic experiments put forward in Rákosi (2012, 2014) and Kertész & Rákosi (2012). This model takes into consideration the current findings of philosophers of science analysing experiments carried out in natural sciences, and seeks analogies with psycholinguistic experiments. In Section 3, the list of well-known criteria put forward by cognitive scientists and

¹ Indeed, it must be mentioned that natural scientists are also liable to the same flaw.

² See also Rákosi (2012: Section 2, 2014: Sections 2-3).

³ For more on this, see Hacking (1983, 1992), Collins (1985), Bogen & Woodward (1988), Pickering (1989),

psycholinguists will be integrated into the metascientific model delineated in Section 2. In Part II of this paper, the proposed system of criteria will be applied to psycholinguistic experiments on metaphor processing conducted between the years 1989 and 2004 in order to exemplify their workability.

2 Experimental methods in psycholinguistic research

2.1 Current views on experiments in science²

The most salient feature of experiments is their *complexity*. They involve a highly complex network of different kinds of activities, physical objects, background knowledge, methods, norms, theories, skills, argumentation processes, interpretative techniques, etc. The *experimental design* is a comprehensive preliminary description of all facets of the process of experimentation which allows for a delineation of the events which are supposed to occur and a rough estimation of the results – that is, it is a special kind of thought experiment. The *experimental procedure* is a material procedure where an experimental apparatus is set up and its operation is monitored and recorded under controlled circumstances. During the experimental procedure *perceptual data* are gained which undergo *authentication* and *interpretation*. Authentication means that the experimenter has to check whether sources of noise do not distort the results: whether they have been ruled out successfully or, if this is not possible, whether their effect can be eliminated with the help of statistical methods. As a result of this process, *experimental data* are obtained which are then confronted with the given theory or with some rival theories.

Nevertheless, it often happens that the interpretation and authentication of the perceptual data indicates shortcomings in the experimental procedure, in the experimental design, or in the theoretical model of the phenomena or of the apparatus. If things do not run smoothly, one turns back to some earlier stage of the experimentation process and modifies a component until there is mutual support among the constituents. All facets of the experiment can be re-examined. Thus, scientific experiments are *cyclic processes*.

A third and often downplayed characteristic of experiments is their *uncertainty*, *fallibility* and *idiosyncrasy*. The continuous emergence of conflicts, inconsistencies and problems, and the process of striving for their elimination is unavoidable with all kinds of research, and experiments are no exception. The conduct of experiments requires practical skills and experience, but also creativity. A series of idiosyncratic decisions has to be made; these decisions cannot always be based on pre-existing methodological rules or reference to precedents, but remain necessarily subjective and arbitrary to a certain extent:

As a knowledge-producing activity, experiment engages the inchoate, the practical, and the particular. The disorderly, inchoate, and personal character of scientific discovery and the complexity of experimental work needed to elicit meaning from phenomenological disorder have persuaded many that there is nothing philosophically interesting to recover [...]. Thus, creative, exploratory, and constructive aspects of experimentation are largely neglected by philosophers of science. (Gooding 2000: 122f.)

Consequently, experiments are *not* completely objective and intersubjectively controllable methods providing perfectly secure and infallible experimental data. Nevertheless, the amount and impact of the potential error sources can be considerably reduced (but not completely

² See also Rákosi (2012: Section 2, 2014: Sections 2-3).

eliminated) by conscious and careful checking and cross-checking.³ To sum up, experiments are interpreted by the literature on experiments in science as *cyclic processes providing uncertain, fallible results*.

Indeed, it is not only the conduct of experiments which requires practical skills, experience, and creativity, but their evaluation, too. Namely, experimental reports are considerably richer than the experimental procedure itself, but, at the same time, they also remain strongly schematic and informationally reduced. What does this mean? On the one hand, the experimental report involves, besides a description of the experiments and its results, *additional elements* such as a short overview of the rival theories, the description of a problem to be solved, and the evaluation and impact of the solution reached, all of which are also achieved with the use of rhetorical tools. On the other hand, it is the experimenter who *selects* the relevant moves and events which are accounted for in the experimental report. Further, she/he has the privilege of deciding what counts as an accidental, contingent and insignificant mistake (which may remain unmentioned) and, what has to be regarded as a possible source of error that has to be reported together with its correction (such as a control experiment), or even as a fatal failure that must lead to the rejection of the data obtained. There are, of course, norms – partly formulated explicitly, partly only implicit – governing experiments as well as experimental reports. The fulfilment of the former, however, cannot be checked directly, but only indirectly, with the help of the latter:

[...] a laboratory notebook and a published journal article are two very different literary forms, serving different purposes and subject to different conventions. The published version should not be viewed simply as a tidied up version of the laboratory notes, since the former contains many conventional elements that would find no place in the latter. The publication is a retrospective narrative, an impersonal, passive reconstruction which draws attention to those theories, tests and data which are considered appropriate for consumption by the scientific community. (Cantor 1989: 160)

2.2 *A model of psycholinguistic experiments*

The basic idea of the model of psycholinguistic experiments as presented in Rákosi (2012, 2014) and Kertész & Rákosi (2012) is that the structure of psycholinguistic experiments basically corresponds to the structure of experiments in the natural sciences, although the role and impact of certain stages diverge.⁴ The model tries to grasp the relationship between the experimental process and the experimental report with the help of *argumentation theoretical tools*. Namely, it assumes that the cyclic process of experimenting as well as the transformation of the experimental procedure into the experimental report involves a *cyclic argumentation process* dealing with different kinds of *uncertain pieces of information* stemming from diverse sources (such as physical activities, events, measurements, theories, statistical tools, etc.). Conversely, the evaluation of experiments might start from the reconstruction of the argumentation presented in the experimental report and a search for traces of the original process of experimenting and the non-public argumentation process related to it. From this it follows that psycholinguistic experiments can be described as *open, cyclic processes*, organised and con-

³ For more on this, see Hacking (1983, 1992), Collins (1985), Bogen & Woodward (1988), Pickering (1989), Bogen (2002), Franklin (2002, 2009), Arabatzis (2008), etc.

⁴ For example, in psycholinguistic experiments, the authentication of perceptual data consists of a checking of the experimental setting (the elaboration and presentation of the stimulus material, the use of fillers, etc.) and only to a lesser degree that of the working of the experimental apparatus. The importance and role of the latter is considerably greater in physics. Rákosi (2012: Section 4) summarises the similarities as well as the differences between experiments in physics and psycholinguistics.

ducted by a *plausible argumentation process*.⁵ This process governs the relationship among hypotheses of the experimental design, the theoretical model of phenomena, the theoretical model of the experimental apparatus, the statements describing the events of the experimental procedure, the statements capturing the results of the interpretation and authentication of perceptual data, as well as the theory being tested and its rivals, etc. See Figure 1.

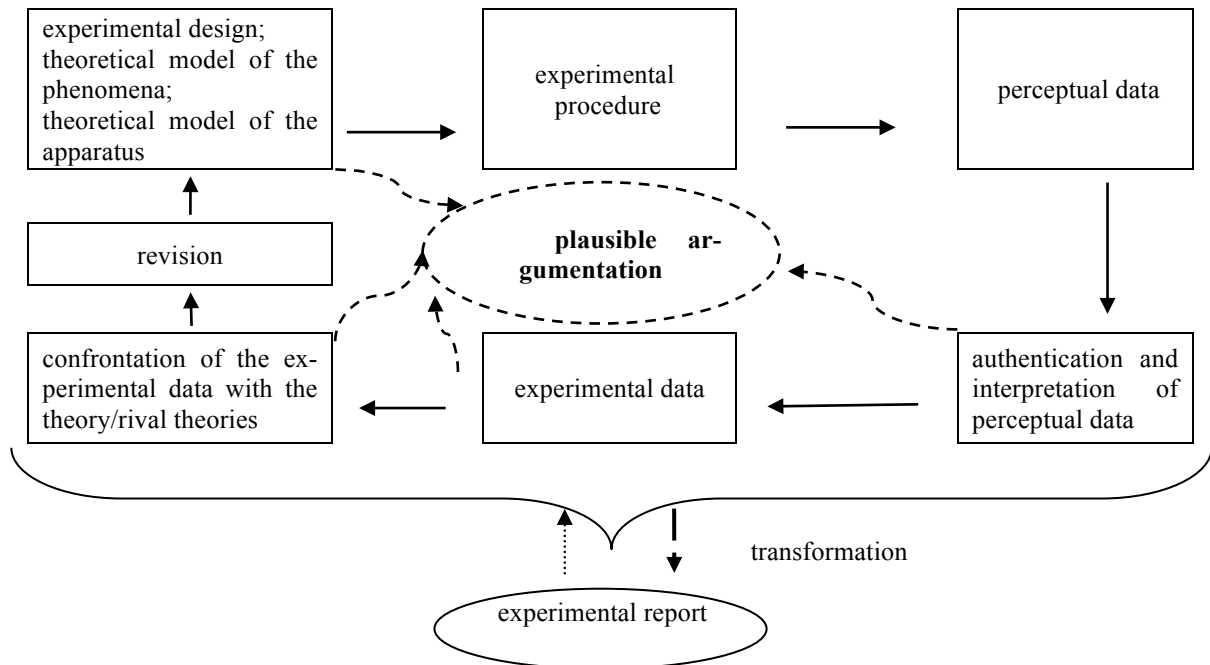


Figure 1: The experimental process⁶

The *p-model* by Kertész & Rákosi (2012, 2014) provides us with tools for representing and comparing the *acceptability* of statements such as previous conjectures, perceptual data, experimental data, hypotheses of linguistic theories, hypotheses about linguistic phenomena, etc. Namely, the *p-model* interprets scientific statements not as propositions but as statements with an informational content and a *plausibility value*. This plausibility value is first determined by the reliability of the source from which it originates and it may change (increase, decrease, or even vanish) during the argumentation process; the plausibility value indicates the level of acceptability of the given statement.⁷ *Plausible argumentation* is a cyclic process that continuously *re-evaluates* the plausibility of statements in the light of new pieces of information and tries to resolve conflicting assessments of the plausibility of hypotheses.⁸

The *effectiveness* and *comprehensiveness* of the plausible argumentation process organising the conduct and control of real experiments largely determines the reliability of the results obtained. This argumentation process is, however, not public but remains the experimenters' *private affair*. It is transformed into an experimental report/paper summarising the results of

⁵ For lack of space, we cannot go into the details of the model of psycholinguistic experiments and the *p-model* of plausible argumentation; we can only delineate their basic ideas here.

⁶ Simple arrows indicate successive stages of the experimental process; dotted arrows signify the non-public argumentation process which organises the experimental process.

⁷ 'Reliability' is used here in a broader sense than in the literature on experiments; that is, the stability of the measurements across conditions is only one of its components.

⁸ That is, it may happen that a statement is plausible on the basis of a source but at the same time, some other source makes its negation plausible.

the experiment and making them available to the scientific community. Clearly, this transformation can be regarded as acceptable if *it does not change the plausibility value of the statements (data, hypotheses) of the original argumentation*. This is of utmost importance because there is a danger that the researcher eliminates relevant information from the published report so that important decisions remain outside public control, and she overestimates the plausibility of the results. From this it follows that the acceptability of the experimental report is also influenced by the *reliability* and *transparency* of the *transformation* of the non-public argumentation process into its public version.

Since the sources available in scientific theorising and experimenting are not completely reliable, errors and inconsistencies in experiments are inevitable and have to be regarded not as fatal failures, but as motivation for their *revision and further development*. Therefore, the active search for possible error sources and inconsistencies is one of the most important driving forces of science, and, in particular, experiments.

3 A system of criteria for the evaluation of psycholinguistic experiments

The key question is, of course, how to decide when an experiment is to some extent reliable as a source and yields plausible (but not certainly true) experimental data and when it is unreliable as a source and is not capable of providing plausible data. A concomitant question is whether the experimental data gained are capable of providing evidence for or against the theory or theories at issue – that is, whether there is a strong enough link between the experimental data and the hypothesis/hypotheses of the theory or rival theories. On the basis of the model presented in Section 2, the evaluation of psycholinguistic experiments involves the following steps.

1) *Reconstruction of the stages of the experimental process in the experimental report*. Although the experimental report can only provide an informationally reduced picture of the experimental process, both the accomplishment of the diverse stages of the experimental process and the cyclic returns conducted by the experimenter in order to eliminate problems revealed should be presented in a detailed enough fashion so that the steps taken can be identified and analysed.

2) *Re-evaluation of the experimental design*. The experimental design should be presented in such a way that the reader is capable of repeating the related thought experiment as well as check validity (such as construct validity, content validity, criterion validity). For example, it should be possible for the reader to check whether the experiment is capable of eliciting participants' natural linguistic behaviour; expectancy effects can be ruled out; semantic priming does not influence participants' performance; participants do not make use of strategic considerations, post-reading checks, or their own implicit theories about the related linguistic phenomena instead of relying on their spontaneous linguistic behaviour, etc.⁹

3) *Re-evaluation of the experimental procedure, the authentication and interpretation of the perceptual data*. The experimental report usually contains hints at revisions of the original experimental design or the experimental procedure. Thus, the evaluation of the experiment has to examine whether possible error sources have been revealed, and whether their impact on the results has been controlled with the help of control experiments or statistical tools. The

⁹ For details, see Kaiser (2013: 139, 141, 143), Haberlandt (1994: 9, 18), Hasson & Giora (2007: 305, 311, 316), Keenan et al. (1990: 384).

interpretation of the perceptual data has to take into consideration, among other things, that there is always only an indirect link between the perceptual data obtained and the linguistic phenomena investigated (such as mental processing of metaphorical expressions). Further, the statistical analysis of the perceptual data is a complex and formidable task with many problematic points, pitfalls and alternatives. Therefore, the conduct of statistical control analyses, alternative analyses and meta-analyses is vital. A further important point is checking the reliability (generalizability) of the results.¹⁰

4) *Re-evaluation of the plausibility of the experimental data and their confrontation with the theory/rival theories.* Since experiments are not completely reliable data sources, they may produce only plausible results. The strength of the support or counter-evidence they may provide to a hypothesis/theory depends on two things: the plausibility of the experimental datum itself, and the strength of the link between the hypothesis/theory and the experimental data. Thus, for example, it has to be checked whether the plausibility value of the experimental data and other data/hypotheses made use of in the experiment is not overestimated in the experimental report; the experimental data (which result from and are bound to a certain situation) can be generalised; alternative explanations can be ruled out (so that the experimental data support only one of the rival hypotheses/theories), etc.

5) *Proposals for the continuation of the experimental process by new cycles.* Since the p-model interprets experiments as open and cyclic processes, the analysis and evaluation of experiments is nothing other than the continuation of the experimental process by new argumentation cycles, and, if possible, the elaboration of proposals for the continuation of the experimental process. Thus, the core of the analysis and evaluation of psycholinguistic experiments are *thought experiments*: one tries to imagine whether and how the experiments described in the experimental report took place and what might have happened, whether there might have been problems which could have distorted the results, etc.

6) *Conduct of replications or modified versions of the experiment.* Thought experiments are, of course, fallible and have their limitations. Thus, while in certain cases such analyses may provide relatively strong counter-arguments (but no ultimate refutations!), which seriously question the reliability of the experiment at issue, in other cases they only indicate weak points and suggest a control experiment or some other kind of revision. Similarly, post hoc statistical analyses of the experimental data are not decisive but have to be taken seriously. Consequently, it might be necessary to transform these thought experiments into real experiments: into a repetition of the original experiment or into a revised version of the experiment, and then compare their outcomes. This means that linguists should not only make their experiments replicable, but that *actual replications* are needed either in an unaltered form or following modifications of the original experimental design.

7) *Comparison of the experimental data with the results of earlier experiments.* Experimental data originating from different experiments cannot be compared mechanically but statistical meta-analyses have to be performed.

To sum up, the evaluation of the weight, impact and treatment of the problematic points of psycholinguistic experiments requires the analysis and weighing up of all details of the given experiment. There are minor flaws that merely *decrease the plausibility* of the affected experimental data, while there are other errors that have to be deemed serious faults that question the usability of the data gained or even *make the experiment unreliable as a data source*.

¹⁰ I use the term ‘reliability’ here in the traditional, narrower sense – that is, it refers to the generalizability of the results to other situations.

Thus, the evaluation of the experiment can and should be accomplished in such a way that not only is its reliability as a data source judged but possible improvements are proposed which, during further cycles, may lead to the continuation and re-evaluation of the experimental process and result in (more) plausible experimental data. See Figure 2.

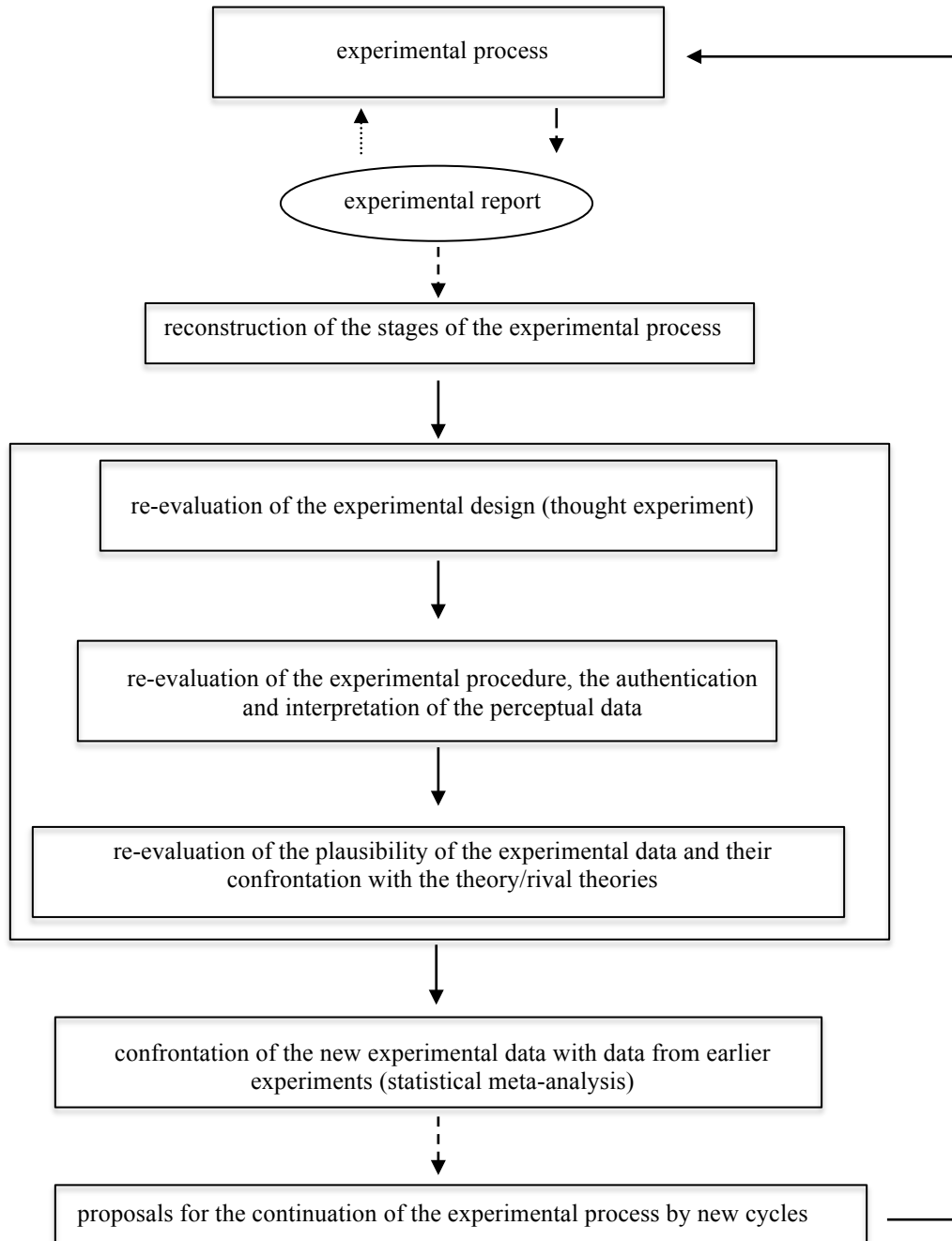


Figure 2: The evaluation of psycholinguistic experiments¹¹

¹¹ Simple arrows indicate successive stages of the experimental process; dotted arrows signify the non-public argumentation process which organises the experimental process.

One might raise the objection that some proposed steps do not provide radically new criteria but rather summarise well-known requirements. Clearly, the collection and systematization of well-established methodological rules, fruitful practices, insights from the philosophy of science, experiences of scientists working in other fields of research, etc. is inevitable, but clearly not sufficient. What is needed, is to *put them to work*.

Nevertheless, it is important to bear in mind that the model of psycholinguistic experiments presented in Section 2.2 and the criteria of evaluation based on it are not a “Wunderwaffe” solving all problems of linguistic experimenting. Therefore, their application does not provide general methodological rules which could be used in every situation, must not be violated, and would guarantee flawless and totally reliable results. Indeed, although experiments are fallible and can provide only plausible experimental data, this does not mean that the above criteria can be violated without consequences. All possible error sources and problems have to be revealed and examined as thoroughly as possible; no weak point and no infringement of the norms should be concealed or ignored. This does not mean that experiments burdened with problems should be immediately rejected; they have to be given appropriate attention and their possible solutions have to be elaborated and compared – or, if this is not possible on the basis of the information at our disposal, this finding has to be declared.

The task of Part II will consist of showing the workability of these ideas with the help of the evaluation of psycholinguistic experiments on metaphor processing conducted between 1989 and 2005. Since Rákosi (manuscript) and Rákosi (in preparation) deal with the theoretical and practical problems related to replications of psycholinguistic experiments, only those experiments will be analysed for which no replication is available yet.

Acknowledgements

Work on this paper was supported by the MTA-DE Research Group for Theoretical Linguistics and the János Bolyai Research Scholarship of the Hungarian Academy of Sciences.

References

- Arabatzis, T. (2008): Experiment. In: Psillos, S. & Curd, M. (eds.): *The Routledge companion to philosophy of science*. Routledge, London & New York, 159-170.
- Bogen, J. (2002): Experiment and observation. In: Machamer, P. & Silberstein, M. (eds.): *The Blackwell guide to the philosophy of science*. Blackwell, Malden & Oxford, 128-148.
- Bogen, J. & Woodward, J. (1988): Saving the phenomena. *The Philosophical Review* 97(3), 303-352.
- Bowdle, B.F. & Gentner, D. (1999): Metaphor comprehension: From comparison to categorization. In: *Proceedings of the Twenty-First Annual Conference of the Cognitive Science Society*, 90-95.
- Cantor, G. (1989): The rhetoric of experiment. In: Gooding, D., Pinch, T. & Schaffer, S. (eds.): *The Uses of Experiment*. Cambridge: Cambridge University Press, 159-180.
- Collins, H.M. (1985): *Changing order: Replication and induction in scientific practice*. Sage, Beverly Hills & London.
- Franklin, A. (2002): *Selectivity and discord. Two problems of experiment*. University of Pittsburgh Press, Pittsburgh.

- Franklin, A. (2009): Experiments. *Stanford Encyclopedia of Philosophy*.
<http://plato.stanford.edu/entries/physics-experiment/>.
- Geeraerts, D. (2006): Methodology in cognitive linguistics. In: Kristiansen, G., Achard, M., Dirven, R. & de Mendoza Ibáñez, F.J.R. (eds.) (2006): *Cognitive Linguistics: Current Applications and Future Perspectives*. Berlin & New York: de Gruyter, 21-49.
- Gibbs, R.W. (2013): The real complexities of psycholinguistic research on metaphor. *Language Sciences* 40, 45-52.
- Gooding, D.C. (2000): Experiment. In: Newton-Smith, W.H. (ed.): *A Companion to the Philosophy of Science*. Malden & Oxford: Blackwell, 117-126.
- Haberlandt, K. (1994): Methods in reading research. In: Gernsbacher, M.A. (ed.): *Handbook of psycholinguistics*. Madison, Wisconsin: Academic Press, 1-31.
- Hacking, I. (1983): *Representing and intervening*. Cambridge University Press, Cambridge.
- Hacking, I. (1992): The Self-Vindication of the Laboratory Sciences. In: Pickering, A. (ed.): *Science as Practice and Culture*. University of Chicago Press: Chicago, 29-64.
- Hasson, U. & Giora, R. (2007): Experimental methods for studying the mental representation of language. In: Gonzalez-Marquez, M., Mittelberg, I., Coulson, S. & Spivey, M. J. (eds.): *Methods in Cognitive Linguistics*. Benjamins, 304-324.
- Kaiser, E. (2013): Experimental paradigms in psycholinguistics. In: Podesva, R.J. & Sharma, D. (eds.): *Research Methods in Linguistics*. Cambridge: Cambridge University Press, 135-168.
- Keenan, J.M., Potts, G.R., Golding, J.M. & Jennings, T.M. (1990): Which elaborative inferences are drawn during reading? A question of methodologies. In: Balota, D.A., Flores d'Archaïs, G.B. & Rayner, K. (eds.): *Comprehension processes in reading*. Hillsdale: Erlbaum, 377-402.
- Kertész, A. & Rákosi, Cs. (2012): *Data and Evidence in Linguistics: A Plausible Argumentation Model*. Cambridge: Cambridge University Press.
- Kertész, A. & Rákosi, Cs. (2014): The p-model of data and evidence in linguistics. In: Kertész, A. & Rákosi, Cs. (eds.): *The Evidential Basis of Linguistic Argumentation*. Amsterdam & Philadelphia: John Benjamins, 15-48.
- Pickering, A. (1989): Living in the material world: On realism and experimental practice. In: Gooding, D., Pinch, T., Schaffer, S. (eds.): *The uses of experiment. Studies in the natural sciences*. Cambridge University Press, Cambridge, 275-297.
- Rákosi, Cs. (2012): The fabulous engine: strengths and flaws of psycholinguistic experiments. *Language Sciences* 34, 682-701.
- Rákosi, Cs. (2014): On the rhetoricity of psycholinguistic experiments. *Argumentum* 10, 533-547. http://argumentum.unideb.hu/2014-anyagok/angol_kotet/rakosicsi.pdf.
- Rákosi, Cs. (manuscript): *'Experimental complexes' in psycholinguistic research on metaphor processing*.
- Rákosi, Cs. (in preparation): *Replication of psycholinguistic experiments and the resolution of inconsistencies*.