

1 Assessing Therapy Effects in Clinical Trials Using  
2 a Measure Theoretic Quanta Analysis

3  
4 Lem Moyé, MD, PhD  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35

36 **Correspondence to:**  
37 Lem Moyé, MD, PhD  
38 UTHealth School of Public Health  
39 1200 Pressler St.  
40 Houston, Texas 77030  
41 Phone: (713) 500-9518  
42 Email: [Lemmoye@msn.com](mailto:Lemmoye@msn.com)  
43

44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61

## Abstract

Statistical hypothesis testing is a fixture in clinical trials. However, its continued application has produced an artificially constrained clinical trial analysis paradigm that is tightly bound to type I error management. Not only is this difficult to helpfully apply in modern complex clinical trials, but it is not responsive to the *primum movens* of the investigators.

A measure theoretic approach is developed here that is based on the principals of health care research analysis and not those of statistical hypothesis testing theory. This new rubric permits all data collected by the trial that is responsive to a specific scientific question to quantitatively contribute to that question's answer. Thus, estimates for each of the following are obtainable: 1) the total available evidence in a clinical trial to answer the question, 2) the strength of that evidence, 3) the strength of evidence that supports benefit and the strength of evidence supporting harm, and 4) the magnitude of the beneficial effect and the magnitude of harm. The incorporation of sampling error in these estimates is achieved without formal hypothesis testing, obviating the need for type I error consideration with its attendant multiplicity corrections.

Keywords. Clinical trial, measure theoretic, benefit-risk ratio

62 **Introduction**

63 Ninety-two years have passed since the writings of Ronald Fisher introduced inference testing to  
64 the applied statistical community [1,2]. This theory of statistical hypothesis testing generated the  
65  $p$ -value that subsequently garnered the support of US Food and Drug Administration (FDA)  
66 regulators, National Institutes of Health administrators, and medical journal editors in assessing  
67 clinical research [3]. With the advent of clinical trials, statistical hypothesis testing became a  
68 fixture among medical researchers and, despite the concerns voiced principally by  
69 epidemiologists [4,5,6,7,8], statistical hypothesis testing remains a fixture of clinical  
70 investigation today, including cardiology, the focus of this manuscript.

71         These statistical hypothesis testing requirement and its focus on  $p$ -values generated a  
72 collection of interpretative conundrums [9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22] for  
73 the cardiology research community. From this turmoil arose additional design and analysis  
74 requisites that a clinical trial must satisfy [23, 24]. These requirements of 1) differentiating  
75 prospectively declared analyses from *post hoc* (exploratory) endpoints and 2) conserving the  
76 overall type I error among a small number of prospectively declared endpoints (multiplicity  
77 corrections) instilled important and necessary discipline in conducting and interpreting clinical  
78 trial results. However, an unfortunate consequence of these tenets is that commonly only a  
79 fraction of the data that are collected in a clinical trial are actually used to directly answer the  
80 study question, a restriction that is required to control type I error propagation [25]. Thus,  
81 although many analyses are conducted, only a small subset of them (and commonly only one of  
82 them) is identified as a “primary”. These restrictions are applied not just in academic research  
83 but in clinical trials conducted by the private sector which follows the contemporary guidance of  
84 the federal Food and Drug Administration [26].

85           This analysis parsimony – a consequence of type I error control – is emblematic of the  
86 fundamental tension between biostatistics and clinical researchers; the inability of statistical  
87 hypothesis testing to directly address the primary, probing question that motivated investigators  
88 to execute their research.

89           The principal inquiry of interest to investigators in clinical trials is “Are participants  
90 better or worse off after exposure to the intervention when compared to the control group  
91 experience?” This is a global question that requires a comprehensive review of all analyses that  
92 bear on this query. However, in cardiology, the combination of 1) the daunting universe of  
93 possible assessments (e.g., heart function, renal function, the ability of the individual to exercise,  
94 how long the individual survives, number of hospitalizations, patient quality of life surveys) and  
95 2) the different types of statistical estimators implemented to assess the exposures effects (e.g.,  
96 Bayes procedures, non-parametric U statistics, regression analyses, imputation evaluations,  
97 survival analyses) have complicated all attempts to provide an answer.

98           The traditional approach of biostatistics is to require the investigator to select one or a  
99 small number endpoints, and then execute statistical hypothesis testing on each, converting every  
100 one of them into the familiar dichotomous decision framework (rejection or non-rejection of the  
101 null hypothesis); type I error is calculated and accumulated as each of these endpoints is  
102 assessed. Thus, to the investigator who has accepted the task of interpreting a complex study of  
103 a complicated disease, statisticians deliver 1) the results of a small number of analyses which the  
104 statisticians believe are dispositive, and 2) an accumulated type I error rate.

105           Although the investigators’ point of view is quantitative, it is also contradistinctive.  
106 Investigators believe that an analysis finding (e.g., the mean difference in the change in exercise  
107 tolerance between the exposed and control group), because of sampling error and other sources

108 of imprecision, can support both a degree of benefit and a degree of harm simultaneously. It is  
109 this dualism – not hypothesis testing dichotomy – to which the investigators resonate, and it is  
110 this dualism – not type I error – that should be accumulated so that investigators can assimilate  
111 their research results. Unfortunately, the standard statistical analysis disappoints the investigators  
112 who find themselves left with 1) no quantitative answer to their principal question and 2) an  
113 accumulated type I error rate in which they have no direct interest. This is the disconnect  
114 between these two scientific disciplines.

115 This paper establishes a rubric which relies on measure theoretic tools to develop  
116 Lebesgue-Stieltjes functions that assess the evidence from all analyses in a clinical trial that are  
117 responsive to a global question, and that provide the degree to which that evidence supports  
118 benefit and harm. This approach provides a direct answer to the investigator’s principal question  
119 with no reliance on statistical hypothesis testing.

120

## 121 **Methods**

122 This development assumes that there is one clinical trial that has been well designed and  
123 concordantly executed (i.e., carried out in accordance with the prospectively written protocol). It  
124 is also assumed that the investigators designed the study to answer one overall question  $q$ , e.g.,  
125 “Does mesenchymal cell therapy improve the well-being of patients with heart failure?” This  
126 manuscript develops answers to the four following inquiries related to question  $q$ :

127

- 128 1. Within the scope of all of the clinical trial’s analyses, what is the content of evidence  
129 that addresses specific question  $q$ ?

- 130           2. Within the scope of all of the clinical trial's analyses, what is the strength of  
131           evidence that actually addresses the specific question  $q$ ?
- 132           3. What is the strength of evidence supporting an affirmative answer to question  $q$   
133           (i.e., a beneficial effect of cell therapy) and the strength of evidence in the trial  
134           suggesting the reverse (i.e., a harmful effect of cell therapy)?
- 135           4. From the evidence of all analyses, what is an estimator of that benefit? What is the  
136           estimate of harm?

137

138           The goal is to create a sample space  $\Omega$  of clinical trial analyses, from which a standard  $\sigma$ -  
139 algebra  $\Sigma$  is formulated. With this as a foundation, a formal measure  $\psi$  is developed on which  
140 analysis-measurable functions operate and can be integrated with respect to  $(\Omega, \Sigma)$ . Their  
141 integrals produce answers to queries 1-4.

142           The full development is available (Appendix). To recapitulate, since a clinical trial's  
143 product is a collection of analyses,  $\{\omega_i\}$ , a sample space  $\Omega$  containing all of the analyses is  
144 generated by the study. Each element  $\omega_i \subset \Omega$  contains the constitutive components of the  $i^{\text{th}}$   
145 analysis. One element of  $\omega_i$  is the question that motivated the analysis. Denote this element of  
146  $\omega_i$  as  $q_i$ . Another collection of components of  $\omega_i$  is the group of analysis characteristics  
147 denoted as  $\delta_i(j)$ ,  $j = 1, \dots$ , where  $j$  indexed the design and operational features of the analysis.  
148 There are many of these characteristics of the analysis, e.g., planning of the analysis (prospective  
149 versus retrospective), and the type of analysis (e.g., survival analysis, mean different analysis,  
150 etc., subgroup analysis). The remaining components of  $\omega_i$  are the participants used in the  
151 analysis, the variables used in the analysis (not their values, but their identities), and the analysis'  
152 estimate of effect size and its standard error.

153           From this perspective one can, for example, collect a set of analyses  $A$  containing all  
154 subgroup assessments evaluating the role of an antidiabetic medication on changes in micro  
155 albuminuria, or a set of analyses  $B$  containing all analyses conducted that examined the  
156 difference in the change over time in systolic blood pressure by therapy group. This

157 multicomponent structure of  $\omega_i$  offers a wide latitude in the creation of sets of analyses. These  
 158 analysis collections, or “regions of analysis” can then serve as the domain on which Lebesgue-  
 159 Stieltjes integrals operate.

160 With this framework, define the content of the analysis  $\omega_i$ , as  $\psi(\omega_i)$ , and write

161 
$$\psi(\omega_i) = n_i v_i.$$

162 This defines the content of an analysis as the product of the number of participants whose data  
 163 contribute to the evaluation multiplied by the number of variables that are required for the  
 164 analysis. Denote the content of this intersection as  $\psi(\omega_i \cap \omega_j)$  where  $\omega_i$  and  $\omega_j$  are not disjoint  
 165 and define

166 
$$\psi(\omega_i \cap \omega_j) = n_{ij} v_{ij}.$$

167 Here,  $n_{ij}$  is the number of participants and  $v_{ij}$  the number of variables common to both analyses.

168 In general, the content of the intersection of  $k$  analyses  $\omega_1, \omega_2, \omega_3, \dots, \omega_k$  is

169 
$$\psi\left(\bigcap_{i=1}^k \omega_i\right) = n_{i\dots k} v_{i\dots k}$$

170 It has been demonstrated that  $\psi(\omega_i)$  meets the formal definition of a measure (Appendix).

171 Since the measure of an analysis is simply based on the number of observation and  
 172 variables it contains, it is easily anticipated that analyses are in general not pairwise disjoint,

173 complicating the computation of  $\psi\left(\bigcup_{k=1}^n \omega_k\right)$ . To expedite this computation, define  $\{B_k\}$  as the

174 sequence of sets created from the increasing sequence of sets  $C_k = \bigcup_{i=1}^k \omega_i$  where  $B_k = C_k \cap C_{k-1}^c$ .

175 Then while it is true that  $\bigcup_{k=1}^n \omega_k = \bigcup_{k=1}^n B_k$ , it is also true that  $\psi\left(\bigcup_{k=1}^n \omega_k\right) = \sum_{k=1}^n \psi(B_k)$  since  $\{B_k\}$

176 consists of pairwise disjoint sets. (Figure 1.) This collection of sets  $\{B_k\}$  represent the analysis

177 fragments or quanta that make separate contributions to the measure of the union of all analyses

178 responsive to a question  $q$ . Furthermore, the measure of any analysis quanta  $B_k$  can be computed

179 as

180 
$$\psi(B_k) = (nv)_k - \sum_{j_1=1}^{k-1} (nv)_{j_1 k} + \sum_{j_1=1}^{j_2-1} \sum_{j_2=2}^{k-1} (nv)_{j_1 j_2 k} - \sum_{j_1=1}^{j_2-1} \sum_{j_2=2}^{j_3-1} \sum_{j_3=3}^{k-1} (nv)_{j_1 j_2 j_3 k} + \dots$$

181 Where  $(nv)_{ij\dots}$  is simplifying notation for  $n_{ij\dots} v_{ij\dots}$  (Appendix).

182 Thus the measure or content of the collection of non-disjoint analyses,  $A = \bigcup_{i=1}^k \omega_i$  can be

183 assembled from the sum of measures of mutually disjoint combinations of analysis quanta

184  $\{B_i\}, i = 1, 2, 3, \dots, k$ , thereby permitting the expression of  $\psi(A) = \int_A d\psi = \int_{\bigcup_{i=1}^k B_i} d\psi = \sum_{i=1}^k \psi(B_i)$

185 where  $\bigcup_{i=1}^k B_i = A$ . An adaptation of this measure for the circumstance in which the variables both

186 within and across analyses are correlated is available (Appendix).

187

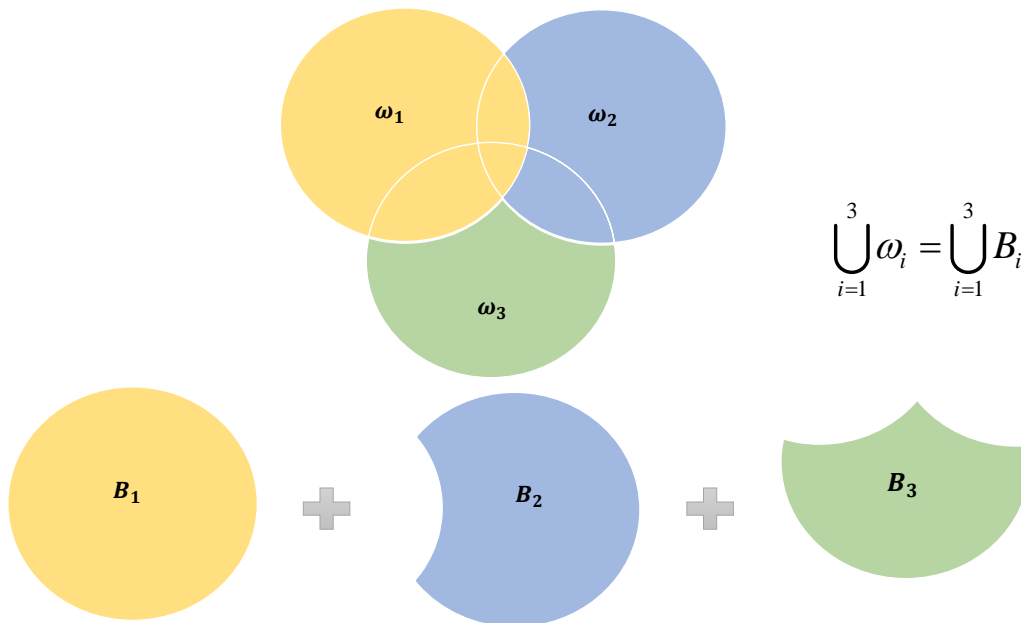


Figure 1. Decomposing the overlapping analyses  $\{\omega_1, \omega_2, \omega_3\}$  into non-overlapping analysis components  $\{B_1, B_2, B_3\}$

188  
189

190 With this as background,  $\psi$  – measureable functions will now be created to address the four  
191 inquiries.

192

193 *Inquiry 1. What is the content of evidence in the trial that addresses question q?*



194 A clinical randomized trial designed to address the question “Does mesenchymal cell therapy  
 195 delivered intravenously improve the cachexic state of advanced heart failure?” conducts many  
 196 analyses to address this inquiry. However, although the analyses are germane, their different  
 197 methodologies (imputation analyses, regression analyses, Bayes analyses, etc.) challenge any  
 198 attempt to combine them. The ultimate goal is to use Lebesgue-Stieltjes integrals to accumulate  
 199 these findings. As a preamble to these steps, from the set of analyses  $\{\omega_i | q_i \subset q\}$  assemble the  
 200 collection of analyses quanta  $\{B_i\}$ , and compute the measure of these quanta using  $\psi$  – measure  
 201 (Appendix).

202 Begin with the identification and collection of the subset of all analyses conducted that  
 203 address question  $q$ . Denote this subset as  $A_q = \{\omega_i / q_i = q\}$ . The content of evidence that  
 204 addresses question  $q$  is obtained by measuring or accumulating the content of analyses that  
 205 contribute to  $A_q$ . Thus, if the content of evidence that addresses question  $q$  is  $\Gamma_q$  then write  
 206  $\Gamma_q = \int_{A_q} d\psi = \psi(A_q)$  where the integral signifies Lebesgue-Stieltjes integration over the collection  
 207 of analyses  $\omega_i \subset A_q$ .

208 Since  $\psi(A_q)$  involves the measure of overlapping sets  $\omega_i \subset A_q$  the inequality  
 209  $\int_{A_q} d\psi = \psi(A_q) = \psi\left(\bigcup_{i=1}^n \omega_i \mathbf{1}_{\omega_i \subset A_q}\right) \leq \sum_{i=1}^n \psi(\omega_i) \mathbf{1}_{\omega_i \subset A_q} = \sum_{i=1}^n n_i v_i \mathbf{1}_{\omega_i \subset A_q}$  (where there are  $n$  member analyses  
 210 contained in  $A_q$ ) is available but not sharp. However the creation of the collection of disjoint  
 211 analyses quanta  $\{B_i\}$  of the previous section permits

212  $\int_{A_q} d\psi = \psi(A_q) = \psi\left(\bigcup_{i=1}^n \omega_i \mathbf{1}_{\omega_i \subset A_q}\right) = \sum_{i=1}^n \psi(B_i) \mathbf{1}_{\omega_i \subset A_q}$ . Thus, the measure of evidence that addresses

213 question  $q$  is  $\Gamma_q = \int_{A_q} d\psi = \psi(A_q) = \sum_{i=1}^n \psi(B_i) \mathbf{1}_{\omega_i \subset A_q}$  with analysis  $\omega_i$  is represented by its

214 quantum  $B_i$ .

215

216 *Inquiry 2. What is the strength of evidence for any analysis in the trial that addresses question*

217 *q?*

218 The measure  $\psi$  is based only on the number of participants and the number of variables

219 that is contained in an analysis, representing only the data that is incorporated in a body of

220 analyses  $A_q = \{\omega_i / q_i = q\}$ . Inquiry 2 addresses not just evidence brought to bear to address

221 question  $q$  but the strength of that evidence provided by the set of analyses  $A_q = \{\omega_i / q_i = q\}$ .

222 Begin by asserting that the strength of evidence contained in each analysis is determined by the

223 research community and transmitted to  $A_q = \{\omega_i / q_i = q\}$  through the decisions of the

224 investigators. For example, in a clinical trial, prospectively declared analyses are commonly held

225 to be of greater value than *post hoc* or exploratory analyses. As another example, measures of

226 organism function (e.g., survival, amputation free survival, walking distance, quality of life), can

227 be of greater value than changes in measures of organ function (e.g. left ventricular ejection

228 fraction (LVEF)), which are themselves of greater value than isolated findings for biomarkers.

229 Thus the strength of evidence contained in an analysis is determined *a priori* by the investigators.

230 This in turn determines the importance of the contribution that the analysis makes to answering

231 the scientific question  $q$ . The nomenclature commonly used to communicate this concept is the

232 use of adjectives such as “prospective”, “primary”, “secondary”, etc.

233 Thus, the formal evaluation process in a clinical trial involves a priority ordering of the

234 analyses’ contributions from the most influential and important to the least contributory. This

235 clinical trial methodology is incorporated here [Appendix]. The choice of a well defended  
 236 analysis priority *a priori* is equivalent to creating a function  $T$  that oversees the reordering of the  
 237 set of analyses  $\{\omega_i\}$  from essentially a random sequence of analyses to a specifically ordered set,  
 238 i.e.,  $T(\omega_1, \omega_2, \omega_3, \dots, \omega_n) = \omega_{[1]}, \omega_{[2]}, \omega_{[3]}, \dots, \omega_{[n]}$  where the subscript  $[i]$  denotes the  $i^{\text{th}}$  analysis in  
 239 the priority order from highest to lowest priority. Note that the function  $T$  also converts the  
 240 sequence  $\{B_i\}, i = 1, 2, 3, \dots$  to  $\{B_{[i]}\}, i = 1, 2, 3, \dots$  the sequence of disjoint analyses quanta  
 241 corresponding to the sequence of analyses ordered by priority. The reordering of  $\{B_i\}$  is critical,  
 242 because an implication of the collection of quanta  $\{B_i\}$  being pairwise disjoint is that their  
 243 contribution to  $\psi\left(\bigcup_{i=1}^n \omega_i\right)$  depends on their location in the priority sequence. Thus the ordered  
 244 sets  $\{B_{[i]}\}, i = 1, 2, 3, \dots$  manifests the *a priori* sense of the importance of the evidence (as reflected  
 245 by the magnitude of  $\psi(B_{[i]})$  to be provided by the analysis.

246 Therefore, the strength of evidence offered by any analysis  $\omega_i = \psi(B_{[i]})$  and the relative  
 247 strength of evidence provided by the  $\omega_i^{\text{th}}$  analysis to address question  $q$  is **RSE** where

$$248 \quad \mathbf{RSE}[\omega_i] = \frac{\psi\left(B_{[i]} \mathbf{1}_{\omega_i \subset A_q}\right)}{\int_{A_q} d\psi} = \frac{\psi\left(B_{[i]} \mathbf{1}_{\omega_i \subset A_q}\right)}{\Gamma_q}.$$

249  
 250 *Inquiry 3. What is the strength of evidence in the trial supporting an affirmative answer to*  
 251 *question  $q$ ? What is the strength of evidence in the trial suggesting a negative answer?*

252

253 The two parts of inquiry 3 will be addressed in turn. Assume that question  $q$  concerns the benefit  
254 or harm of an intervention in a clinical trial, e.g., “Does the provision of mesenchymal cells to  
255 patients with heart failure ameliorate their signs and symptoms when compared to the experience  
256 of controls?” The process to be followed to address this question is to first identify the statistical  
257 estimate of effect from each analysis  $\omega_i \subset A_q = \{\omega_i / q_i = q\}$  and then for each estimator, 1)  
258 consider the distorting role of sampling error and imprecision on this estimate, 2) parse the  
259 resulting region into a region supporting benefit, 3) quantify this region, 4) norm this by the  
260 measure of its quanta, and 5) accumulate this evidence over all  $\omega_i \subset A_q$ .

261 One of the components of each  $\omega_i$  is the effect size produced by the analysis, identified  
262 now as  $e_i$ . This quantity  $e_i$  can be the difference between therapy groups of the mean blood  
263 pressure change over time, or the relative risk of death associated with an intervention.  
264 However, due to sampling variability and the measurement’s relative imprecision, this estimate  
265 of benefit cannot be relied upon in and of itself. The impact of these two distorting effects is to  
266 blur the exact position of the population measure of effect that could be deduced from the value  
267 of the statistical estimator from the sample; variability and imprecision each suggest that both  
268 larger values and smaller values of the estimator are admissible for consideration. This range of  
269 values will be termed the estimator’s region of plausible effects. It is not just the estimator that  
270 provides a sense of the effect of the intervention; it is the estimator’s region of plausible values  
271 that is most informative about the possible effect size that would be seen in the population.

272 The region of plausible effect will always provide values of the effect size that are larger  
273 than the statistical estimator, and others that are smaller. In many cases, the distorting effects of  
274 imprecision and sampling error can actually reverse the direction of effect, signifying that not  
275 benefit, but harm might be produced in the population at large.

276 The observation that the statistical estimator produces a plausible region of effect that  
 277 together and simultaneously supports both larger benefit values and smaller ones (that sometimes  
 278 includes harm) is here termed duality. Estimators refract the data on which they are based into  
 279 both larger and smaller effect sizes including effect sizes that are indicative of harm. It is this  
 280 duality that the functions developed in this section will first segregate and capture (a process  
 281 termed analysis parsing) and then accumulate using  $\psi$  – measure.

282 Define the upper  $e_i^+$  and lower  $e_i^-$  bounds of an interval of plausible effect for the  
 283 analysis as  $\omega_i$ , computing

$$284 \quad \begin{aligned} e_i^+ &= e_i + a_i \\ e_i^- &= e_i - b_i \end{aligned}$$

285 where  $a_i$  and  $b_i$  are constants based on variability and imprecision. Note that this interval need  
 286 not be symmetric around the actual estimator  $e_i$ . The region of plausible effect is signified as  
 287  $[e_i^-, e_i^+]$ .

288 This plausible effect interval can be parsed into two subintervals, one a region of benefit,  
 289 the other of harm. In order to locate these sub-regions, knowledge of the value of the statistical  
 290 estimator's effect that is neutral (i.e., denotes neither benefit nor harm) is required. Define this  
 291 value of neutral effect as  $e_i(0)$ . Similarly, let  $e_i(b)$  and  $e_i(h)$  be the values of the greatest  
 292 possible benefit and the greatest possible harm permitted by the estimator respectively. The  
 293 introduction of  $e_i(b)$  and  $e_i(h)$  is necessary since values of harm need not always be less than  
 294 values of benefit. For example, if the  $i^{\text{th}}$  analysis is a total mortality hazard function analysis,  
 295 then  $e_i = 1$  indicates no effect on the time to death,  $e_i(h) = \infty$ , and  $e_i(b) = 0$ . Alternatively, if  
 296  $\omega_i$  is an evaluation of changes in mean differences where the greater differences are salubrious,

297 then the value of  $e_i = 0$  reflects no mean effect,  $e_i(h) = -\infty$ , and  $e_i(b) = \infty$ . Using this notation,  
 298 then the interval  $[\min(e_i(h), e_i(b)), \max(e_i(h), e_i(b))]$  is the range of possible values of the  
 299 estimate.

300 Consider the case where  $e_i(b) > e_i(h)$ . We now define the plausible benefit interval  $\chi_i^{(b)}$   
 301 as;

$$302 \quad \chi_i^{(b)} = [b_i^-, b_i^+] = [e_i^-, e_i^+] \cap [\min(e_i(0), e_i(b)), \max(e_i(0), e_i(b))] = \mathbf{1}_{[e_i^-, e_i^+]} \mathbf{1}_{[e_i(0), e_i(b)]} = \mathbf{1}_{[b_i^-, b_i^+]}$$

303 This is the portion of the plausible effect size region that supports benefit. For example, larger  
 304 values of left ventricular ejection fraction are considered beneficial *ceteris paribus*; its increases  
 305 are beneficial and its decreases are harmful. Thus, if the plausible effect region for a change in  
 306 left ventricular ejection fraction is  $[-1, 7]$  and the region of these changes that are beneficial is  
 307  $(e_i(0), e_i(b)) = (0, \infty)$ , then  $\chi_k^{(b)} = [-1, 7] \cap (0, \infty) = (0, 7]$  is the plausible benefit region. The

308 plausible region for harm is based on  $(\min(e_i(h), e_i(0)), \max(e_i(h), e_i(0))) = (-\infty, 0)$ , and is

$$309 \quad \chi_i^{(h)} = [h_i^-, h_i^+] = [e_i^-, e_i^+] \cap [e_i(h), e_i(0)] = \mathbf{1}_{[e_i^-, e_i^+]} \mathbf{1}_{[e_i(h), e_i(0)]} = \mathbf{1}_{[h_i^-, h_i^+]}$$

310 which in this example is  $\chi_k^{(h)} = [-1, 7] \cap (-\infty, 0) = (-1, 0]$ .

311 Now define the contribution function

$$312 \quad \mathbf{Y}(\chi_i^{(b)}) = \mathbf{Y}\left(\mathbf{1}_{[b_i^-, b_i^+]}\right) = \frac{1}{(b_i^+ - b_i^-)} \left( \frac{b_i^+ + b_i^-}{2} + b_i^- \right)$$

313 as the unit-less benefit function that maps the interval of plausible benefit to an assessment of the  
 314 level of that benefit.  $\mathbf{Y}(\chi_i^{(b)})$  penalizes the benefit estimate derived from  $\omega_i$  for a wide interval,

315 while amplifying benefit if the minimum value of the plausible region is different than  $e_i(0)$

316 (Figure 2).

317

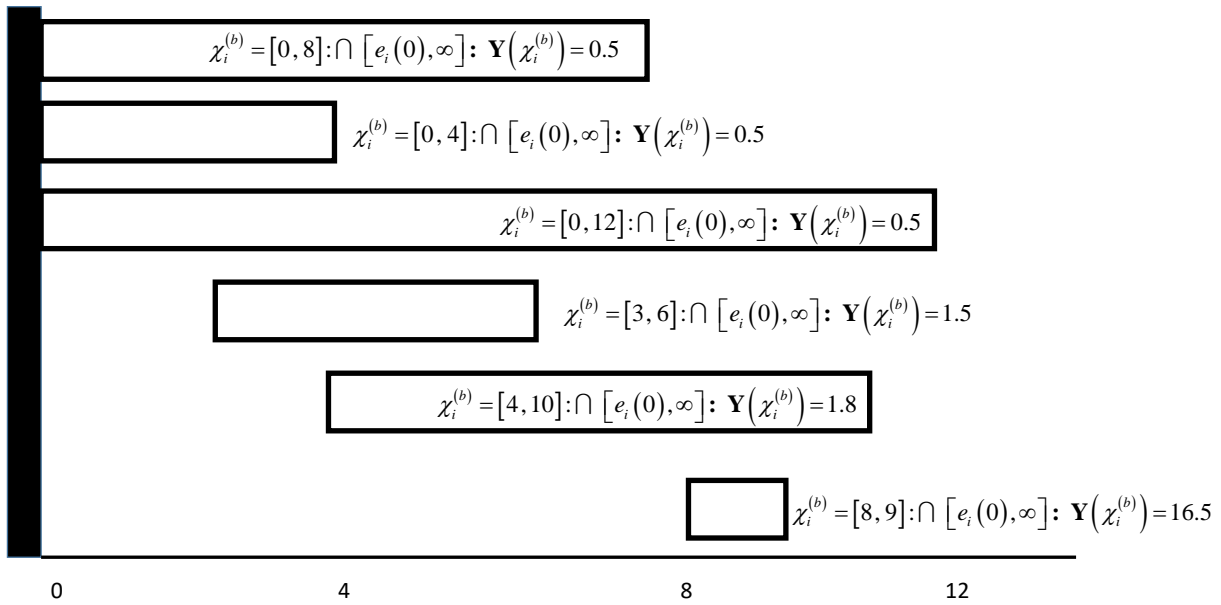


Figure 2. Operation of the benefit function for different levels of analyses effects.

318

319 From Figure 2, the circumstance where  $\chi_i^{(b)} = [b_i^-, b_i^+] = [0, 8]$ ,  $\chi_i^{(b)} = [b_i^-, b_i^+] = [0, 4]$  and

320  $\chi_i^{(b)} = [b_i^-, b_i^+] = [0, 12]$  each generate a contribution function value of only

321  $\mathbf{Y}(\chi_i^{(b)}) = \mathbf{Y}(\mathbf{1}_{[b_i^-, b_i^+]}) = 0.5$ , reflecting some addition of benefit from this region, but penalizing it

322 because their lower bound includes  $e_i(0)$ , the value of no effect. The contribution's function

323 value is greater when  $b_i^- > e_i^0$ , as is the case of the remaining two examples in Figure 2.

324 Analogous contribution computations manage the harm concern.

325 With the benefit interval and contribution function in hand, it remains to compute the

326 benefit over all of the analysis  $\omega_i \subset A_q$ . The integral  $\int_{A_q} \mathbf{Y}(\chi_i^{(b)}) d\psi$  is the assessment of the benefit

327 function on each set  $\omega_i \subset A_q$  with respect to  $\psi$ -measure. A normed version over the measure of  
 328  $A_q$  can be written as  $\mathbf{B}_q = \left[ \psi(A_q) \right]^{-1} \int_{A_q} \mathbf{Y}(\chi_i^{(b)}) d\psi$ . Thus  $\mathbf{B}_q$  is the normed measure of benefit  
 329 derived from all of the analyses responsive to question  $q$ .

330 A similar quantity can be computed to assess harm. With the plausible harm interval

331  $\chi_k^{(h)}$  defined as above where  $\chi_i^{(h)} = [h_i^-, h_i^+] = [e_i^-, e_i^+] \cap [-\infty, e_i(0)] = \mathbf{1}_{[e_i^-, e_i^+]} \mathbf{1}_{[-\infty, e_i(0)]} = \mathbf{1}_{[h_i^-, h_i^+]}$ .

332 define  $\mathbf{Y}(\chi_i^{(h)}) = \left| \frac{1}{(h_i^+ - h_i^-)} \left( \frac{h_i^+ + h_i^-}{2} + h_i^+ \right) \right|$  and  $\mathbf{H}_q = \left[ \psi(A_q) \right]^{-1} \int_{A_q} \mathbf{Y}(\chi_i^{(h)}) d\psi$ . With these

333 quantities, write the benefit to harm ratio as

$$334 \quad \frac{\mathbf{B}_q}{\mathbf{H}_q} = \frac{\left[ \psi(A_q) \right]^{-1} \int_{A_q} \mathbf{Y}(\chi_i^{(b)}) d\psi}{\left[ \psi(A_q) \right]^{-1} \int_{A_q} \mathbf{Y}(\chi_i^{(h)}) d\psi} = \frac{\int_{A_q} \mathbf{Y}(\chi_i^{(b)}) d\psi}{\int_{A_q} \mathbf{Y}(\chi_i^{(h)}) d\psi}$$

335

336 *Inquiry 4. If there is a benefit, what is an estimator of that benefit? If there is harm, what is the*  
 337 *estimate of that harm?*

338 For a query  $q$ , quantitative assessments of the evidence for benefit  $\mathbf{B}_q$  and harm  $\mathbf{H}_q$

339 were produced. Here an actual estimate of the level of benefit (and of harm) will be developed.

340 Recall that the plausible benefit interval  $\chi_i^{(b)}$  is defined as  $\mathbf{1}_{[b_i^-, b_i^+]}$ . There are several

341 functions that provide service in assessing the effect of therapy based on that interval. Let  $\mathbf{I}$  be

342 the condition where an increase in  $e_i$  reflects benefit and  $\mathbf{D}$  reflect the circumstance where a

343 decrease reflects benefit. Then one such function is  $\mathbf{L}_{\max}(\chi_k^{(b)}) = \mathbf{L}_{\inf}(\chi_k^{(b)}) \mathbf{1}_{\mathbf{D}} + \mathbf{L}_{\sup}(\chi_k^{(b)}) \mathbf{1}_{\mathbf{I}}$ . This

344 represents the assessment of greatest benefit from the plausible interval. Alternative, one could



345 conservatively estimate benefit as  $\mathbf{L}_{\min}(\chi_k^{(b)}) = \mathbf{L}_{\sup}(\chi_k^{(b)})\mathbf{1}_D + \mathbf{L}_{\inf}(\chi_k^{(b)})\mathbf{1}_I$ . This represents the  
 346 least effect value for benefit. Choosing the latter for this development, define the estimate of  
 347 benefit from all of the analyses addressing question  $q$  as  $\Lambda_{qB}$

$$348 \quad \Lambda_{qB}(\min) = \left[ \int_{A_q} d\psi \right]^{-1} \int_{A_q} \mathbf{L}_{\min}(\chi_i^b) d\psi.$$

349 This is the accumulation of unit less benefit with respect to the content of each analysis, normed  
 350 by the accumulated content of all analyses.

351 A similar result is obtained for an estimator of harm produced by all analyses

$$352 \quad \{\omega_i \mid \omega_i \subset A_q\}.$$

$$353 \quad \Lambda_{qH}(\max) = \left[ \int_{A_q} d\psi \right]^{-1} \int_{A_q} \mathbf{L}_{\max}(\chi_i^h) d\psi.$$

354 Where  $\mathbf{L}_{\max}(\chi_k^{(h)}) = \mathbf{L}_{\sup}(\chi_k^{(h)})\mathbf{1}_D + \mathbf{L}_{\inf}(\chi_k^{(h)})\mathbf{1}_I$  is the worst case estimate of harm

355 obtained from  $\chi_k^{(h)}$  obtained from the plausible regions of harm. As with the benefit function,  
 356 alternative views of harm are also available.

357

## 358 Discussion

359 This manuscript provides an alternative approach to clinical trial analysis that is based on both  
 360 the principles of measure theory and clinical trial methodology. Its solutions provide answers to  
 361 four inquiries of critical interest to clinical trialists using statistical estimation theory that are  
 362 commonly not quantitatively addressed, while not relying on statistical hypothesis testing.

363           The current clinical trial analysis procedure requires the discrimination of prospective  
364 from exploratory analyses, and control of the familywise type I error among the former. This  
365 commonly restricts the study’s conclusive analyses to a small number of evaluations addressing  
366 precise “primary endpoints” that are pre-designated to represent the principal findings of the  
367 well-designed, concordantly executed clinical trial.

368           Unfortunately, this standard approach is a symptom of the detachment of the goals of the  
369 clinical trial investigators from the work product of contemporary biostatistics. A clinical trial’s  
370 primary analyses address important questions, but their answers are only contributory to the  
371 more global question of “Has the health, well-being, and sense of well-being of participants  
372 improved after exposure to the new therapy when compared to the experience of the control  
373 group?” This is the question that is of greatest interest to research investigators, participants,  
374 health care providers, formulary committees, and the regulatory community. The answer to this  
375 larger question requires a broad and integral appraisal of all responsive analyses.

376           This is not how biostatistics is applied. Its standard approach is to evaluate a small  
377 number of the many components (e.g., survival, or peak walking time, or improvement in  
378 LVEF) of this omnibus question one at a time, converting the question of, for example “What is  
379 the effect of the intervention on survival time?” into a dichotomous question “Is survival  
380 changed by therapy or not?” This bifurcation is modulated by consideration of the confidence  
381 interval, but in the end, it is then assessed using statistical hypothesis testing, whose product is a  
382 “yes-no” answer and a type I error measurement. Thus, in the end, the classic statistical analysis  
383 procedures proffer the combination of 1) a small collection of dichotomous responses to the  
384 primary endpoints and 2) an overall type I error expenditure as dispositive.

385           However, clinical investigators have an abiding interest in neither. Physicians and  
386 researchers understand that due to the role of measurement imprecision and sampling error, an  
387 effect size that is provided by an endpoint in fact stands for not just one value but for a range of  
388 effect sizes. Some of these effect sizes are supportive of benefit, while others – in a different part  
389 of the range – are less supportive and may even be consistent with harm. Thus, an effect size  
390 range can simultaneously contribute to an argument supporting benefit and also a contention for  
391 harm. It is this analysis-generated dualism that researchers require be drawn together into an  
392 ensemble of effects from different analyses that would be responsive to the global question.  
393 From the investigator’s perspective, it is not type I error but the benefit/harm assessment that  
394 requires accumulation across analyses. Biostatisticians provide the former and investigators need  
395 the latter. This is the disconnect that the measure theoretic approach attempts to repair.

396           In this manuscript, the broad concepts of set and measure theory have been contoured to  
397 address clinical trial evaluations. Specifically a measure, and a collection of measurable  
398 functions have been developed with the single goal of incorporating salient features of clinical  
399 trial methodology into the realm of formal mathematical analysis and measure theory. The result  
400 is a system that is mathematically rigorous, flexible, and can be practically applied.

401           The computations involved in this system are straightforward. The steps are as follows:

- 402           1. Identify the set of analyses  $\{\omega_i / \omega_i \in A_q\}$  that address the question  $q$ .
- 403           2. Pre-specify all of the analyses to be conducted and their priority of importance in  
404           addressing question  $q$ .
- 405           3. For the set  $\{\omega_i / \omega_i \in A_q\}$  using compute  $\psi(\omega_i)$ .
- 406           4. Compute the set of quanta  $\{B_{[i]}\}$  and for each member, compute  $\psi(B_{[i]})$ .

407 5. Compute the measure of the body of evidence that addresses question  $q$  as by using

408 the quanta as  $\mathbf{E}_q = \int_{A_q} d\psi = \psi(A_q) = \sum_{i=1}^n \psi(B_i) \mathbf{1}_{\omega_i \subset A_q}$  and the strength of evidence for

409 proffered for each analysis as  $\frac{\psi(B_{[i]} \mathbf{1}_{\omega_i \subset A_q})}{\Gamma_q}$

410 6. Compute the evidence for benefit and the evidence for harm

411  $\mathbf{B}_q = [\psi(A_q)]^{-1} \int_{A_q} \mathbf{Y}(\chi_i^{(b)}) d\psi$ ,  $\mathbf{H}_q = [\psi(A_q)]^{-1} \int_{A_q} \mathbf{Y}(\chi_i^{(h)}) d\psi$ . and the benefit to harm

412 ratio  $\frac{\mathbf{B}_q}{\mathbf{H}_q}$

413 7. Compute the estimates of benefit  $\Lambda_{q\mathbf{B}}$  (min) and harm  $\Lambda_{q\mathbf{H}}$  (max).

414

415 The research community can choose the functional form for assessing the intervals of  
416 benefit and harm.  $\mathbf{Y}(\chi_i^{(b)})$  was specifically chosen here in order to reduce the impact of the  
417 plausible region of benefit by its length, and increase its impact if its lower boundary was  
418 different from that value delineating no effect. There are other choices available though. Triangle  
419 functions and scaled beta functions require attention. However, it is best to keep in mind that this  
420 system is designed to be relatively easy to use; more complicated forms of  $\mathbf{Y}(\chi_i^{(b)})$  increase the  
421 complexity of the evaluations.

422 One challenge in the application of this measure theoretic approach is the interpretation  
423 of the measure of benefit  $\Lambda_{q\mathbf{B}}$  (min) and measure of harm  $\Lambda_{q\mathbf{H}}$  (max). These are derived as unit-  
424 less quantities, but the research community has no experience with their interpretation. It is  
425 therefore proposed that for the immediate future both the traditional analysis and this measure

426 theoretic approach be conducted in the evaluation of clinical trials. This would provide  
427 calibration for the research community as it works to interpret these new values  $\Lambda_{qB}(\min)$  and  
428  $\Lambda_{qH}(\max)$ .

429 Flexibility of analyses is an advantage of this approach. There is no need to focus on a  
430 particular type of estimator. Standardly used estimators, e.g., mean differences, relative risks,  
431 Bayes procedures, regression estimates, imputation generated effects, can each be incorporated.

432 No prior example of the assignment of a formal measure to a clinical analysis (separate  
433 and apart from the equivalence of the Lebesgue and Riemann integral when the Riemann integral  
434 exists) has been identified. While it seems clear that the “amount” of data on which an analysis  
435 relies is an appropriate contributor to the measure of that analysis, defining the measure of  
436  $\psi(\omega_i) = n_i v_i$  is not the only definition available, and there are clearly alternative measures that  
437 one could apply to the  $(\Omega, \Sigma)$  collection of analyses. However, the framework developed here is  
438 simple, reasonable, and produces tractable computations.

439 A requirement of the approach of this manuscript is to ensure that the structure of the  
440 measure theoretic framework be permeable to clinical trial design requirements. The importance  
441 of priority of analysis is critical in research methodology and therefore is incorporated in the  
442 proposed analysis rubric; the size of the independent contribution of analysis  $\omega_i$  depends on  
443 where it lies in the sequence of evaluations. In fact, this measure theoretic approach provides a  
444 mathematical justification for the long established practice of selecting high priority analyses in  
445 clinical trials; these are the analyses which makes the greatest contribution to the  $\psi\left(\bigcup_{i=1}^n \omega_i\right)$ . By  
446 matching the sequence of quanta  $B_{[i]}$  to the priority of evaluations chosen by the investigators,

447 emphasizes on analyses results  $\mathbf{B}_q$ ,  $\mathbf{H}_q$ ,  $\Lambda_{q\mathbf{B}}$  (min), and  $\Lambda_{q\mathbf{H}}$  (max) are placed precisely where the  
448 investigators have *a priori* stipulated. It is recommended to investigators that this sequence be  
449 chosen based on the importance of the analysis in contributing to the understanding of the effect  
450 of the exposure.

451 The investigator determination of analysis priority viewed from a measure theoretic  
452 perspective provides new approaches to challenging problems in trial design. For example, it is  
453 beyond question that safety evaluations in clinical trials are paramount. However, safety  
454 evaluations are not typically part of the type I error control structure in traditionally analyzed  
455 clinical trials; for example, type I error is commonly not first accrued for safety, with the  
456 remainder being distributed across primary endpoints. The safety analysis lies awkwardly outside  
457 the alpha accumulation structure in the traditional paradigm. However, the measure theoretic  
458 structure presented in this manuscript permits the safety evaluation to be prioritized first,  
459 followed by efficacy evaluations. The impact of the efficacy endpoints would be reduced, but  
460 this is wholly consistent with a *primum non nocere* philosophy. In addition any reduced measure  
461 that is seen in the primary efficacy evaluations because of the first consideration of safety is  
462 partially offset by the accumulation of benefit using  $\mathbf{Y}\left(\chi_{t(\omega_t)}^{(b)}\right)$  for the safety evaluation. This  
463 approach avoids the analytic disconnect; the safety evaluations are incorporated mathematically  
464 and smoothly into the scope of the analyses.

465 In addition, when viewed from a measure theoretic perspective, the door is open for the  
466 investigators to examine different scenarios to optimize the size of  $\psi\left(B_{\omega_t}\right)$  for the analyses of  
467 most interest. This optimization requires not just concern for sequencing, but for maximizing or  
468 minimizing the measure of the intersections between the analyses.

469           This paper is not an argument for the abandonment of rigor. The discipline that  
470 epidemiologists and biostatisticians have helped to instill in investigators is laudable; it is not  
471 argued that the stringent execution of a protocol be dismissed. All analyses to be incorporated  
472 should be identified prospectively and thoroughly vetted before the research endeavor  
473 commences. Endpoint measures should be obtained from state of the art equipment known for  
474 their satisfactory precision. If possible, evaluations that would support or refute the purported  
475 mechanism of action should be incorporated. While the need for statistical hypothesis testing  
476 may be removed when one implements these procedures, the rules of epidemiology and the need  
477 for discipline still apply.

478           An advantage of developing a foundation based on clinical trial methodology for the  
479 mathematical interpretation of the trial's results has the advantage of extensibility. The role of  
480 exploratory analyses has been problematic for the traditional analysis rubric, in which  
481 exploratory analyses are not incorporated into the trial's endpoint analysis. In this measure-  
482 theoretic structure, the exploratory analysis  $\omega_i$  can be incorporated into the final result, but its  
483 role in affecting the final result depends on where in the sequence  $[i]$  lies. In addition, the  
484 integration of analyses that appear from the same trial in separate manuscripts can be  
485 accomplished as well, providing an overall picture of the benefit and harm risks posed by the  
486 exposure being studied. Nonrandomized, observational studies in health care are amenable to the  
487 application of this measure theoretic approach as well, although at this stage of development,  
488 neither  $\psi$  – measure nor the integrals  $\mathbf{B}_q$ ,  $\mathbf{H}_q$ ,  $\Lambda_{q\mathbf{B}}(\min)$ , and  $\Lambda_{q\mathbf{H}}(\max)$  explicitly take into  
489 account the universe of biases that can vitiate the results of the observational study. Finally, there  
490 may be meta-analytic implications of this work.

491 Statistical hypothesis testing has played an important role in clinical trials. However, the  
492 connection between its standard application and the goal of clinical trials is broken. The rubric  
493 described here provides a measure theoretic mathematical structure developed specifically for  
494 clinical trials, allowing the investigator access to the global result they require and for which they  
495 designed the study.

496

497

498

499 References

- 
1. Fisher, R A. (1925) *Statistical methods for research workers*. Edinburg. Oliver and Boyd.
  2. Fisher RA. (1926). The arrangement of field experiments. *Journal of the Ministry of Agriculture*. September 503 - 513.
  3. Goodman, S.N. (1999). Toward Evidence-Based Medical Statistics. 1: The  $p$ -value fallacy. *Annals of Internal Medicine* 130:995–1004.
  4. Walker A.M. (1986). Significance tests represent consensus and standard practice (Letter) *American Journal of Public Health*.76:1033. (See also *Journal erratum* 76:1087.
  5. Fleiss J.L. (1986). Significance tests have a role in epidemiologic research; reactions to A.M. Walker. (Different Views) *American Journal of Public Health* 76:559–560.
  6. Fleiss J.L. (1986). Confidence intervals versus significance tests: quantitative interpretation. (Letter) *American Journal of Public Health* 76:587.
  7. Fleiss J.L. Dr. Fleiss response (Letter) (1986). *American Journal of Public Health* 76:1033–1034.
  8. Walker A.M. (1986). Reporting the results of epidemiologic studies. *American Journal of Public Health* 76:556–558.
  9. Multiple risk factor intervention trial. Risk factor changes and mortality results. Multiple Risk Factor Intervention Trial Research Group. *JAMA*. 1982 Sep 24;248(12):1465-77.
  10. Feldman AM, Bristow MR, Parmley WW et.al (1993) Effects of vesnarinone on morbidity and mortality in patients with heart failure. *N Engl J Med* 329:149-55
  11. Cohn J, Goldstein SC, Feenheed S et.al. (1998) A dose dependent increase in mortality seen with vesnarinone among patients with severe heart failure. *N Eng J Med* 339:1810-16
  12. Pitt B, Segal R, Martinez FA. et al. on behalf of the ELITE Study Investigators (1997). Randomized trial of losartan versus captopril in patients over 65 with heart failure. *Lancet* 349:747–52.
  13. Pitt B, Poole-Wilson PA., Segal R, et. al (2000). Effect of losartan compared with captopril on mortality in patients with symptomatic heart failure randomized trial–The losartan heart failure survival study. ELITE II. *Lancet*.355:1582–87.



- 
14. Packer, M., O'Connor, C.M., Ghali, J.K., et al for the Prospective Randomized Amlodipine Survival Evaluation Study Group (1996). Effect of amlodipine on morbidity and mortality in severe chronic heart failure. *New England Journal of Medicine*.335:1107–14.
  15. Packer, M. (2000). Presentation of the results of the Prospective Randomized Amlodipine Survival Evaluation-2 Trial (PRAISE-2) at the American College of Cardiology Scientific Sessions, Anaheim, CA, March 15, 2000.
  16. Packer M., Bristow MR. Cohn JN. et al. (1996). The effect of carvedilol on morbidity and mortality in patients with chronic heart failure. *New England Journal of Medicine*. 334:1349-55.
  17. Moyé LA, Abernethy D. (1996) Carvedilol in Patients with Chronic Heart Failure (Letter). *New England Journal of Medicine*. 335: 1318-1319.
  18. Packer M, Cohn JN., Ccolucci WS. Response to Moyé and Abernethy. *New England Journal of Medicine* 335:1318-1319.
  19. Fisher L. (1999). Carvedilol and the FDA approval process: the FDA paradigm and reflections upon hypotheses testing. *Controlled Clinical Trials* 20:16-39.
  20. Fisher LD, Moyé LA. (1999) Carvedilol and the Food and Drug Administration Approval Process: An Introduction. *Controlled Clinical Trials*. 20:1-15.
  21. Moyé L.A. (1999) P Value Interpretation in Clinical Trials. The Case for Discipline. *Controlled Clinical Trials* 20:40-49.
  22. Fisher LD. Carvedilol and the Food and Drug Administration-Approval Process: A Brief Response to Professor Moyé's article. *Controlled Clinical Trials*. 20:50-51.
  23. Pfeffer M. A Second Prospective Randomized Amlodipine Survival Evaluation (PRAISE-2). *Cardiology Update*. te.
  24. Pfeffer MA, Skali H. PRAISE (prospective randomized amlodipine survival evaluation) and criticism .*JACC Heart Fail*. 2013 Aug;1(4):315-7. doi: 10.1016/j.jchf.2013.05.005. Epub 2013 Aug 5. No abstract available. PMID: 24621934
  - 25 . Hare JM, Bolli R, Cooke JP, Gordon DJ, Henry TD, Perin EC, March KL, Michael P. Murphy MP, Pepine CJ, Simari RD, Skarlatos SI, Szady A, Taylor DA, Traverse J, Willerson JT, Vojvodic RW, Yang PC, Moyé L for the Cardiovascular Cell Therapy Research Network (CCTRN). Phase II Clinical Research Design in Cardiology Learning the Right Lessons Too Well: Observations and Recommendations from the Cardiovascular Cell Therapy Research Network (CCTRN). *Circulation* 2013;127:1630-35.
  - 26 Guidance for industry. Cellular Therapy for Cardiac Disease.  
<http://www.fda.gov/BiologicsBloodVaccines/GuidanceComplianceRegulatoryInformation/Guidances/default.htm>.