



4 October 2022  
EMADOC-1700519818-946771  
Committee for Medicinal Products for Human Use (CHMP)

## DRAFT Qualification opinion for the iBox Scoring System as a secondary efficacy endpoint in clinical trials investigating novel immunosuppressive medicines in kidney transplant patients

Draft agreed by Scientific Advice Working Party (SAWP)	1 September 2022
Adopted by CHMP for release for consultation	15 September 2022
Start of public consultation	6 October 2022
End of consultation (deadlines for comments)	17 November 2022

Comments should be provided using this [template](#). The completed comments form should be sent to [ScientificAdvice@ema.europa.eu](mailto:ScientificAdvice@ema.europa.eu)

<b>Keywords</b>	Qualification of Novel Methodology, iBox, composite biomarker panel, eGFR, proteinuria, renal allograft biopsy, DSA, time-post-transplant, secondary efficacy endpoint, kidney transplant clinical trials, immunosuppressive medicines, allograft failure
-----------------	---

<sup>1</sup> Last day of relevant Committee meeting.

<sup>2</sup> Date of publication on the EMA public website

<sup>3</sup> Last day of the month concerned.



1 **Table of contents**

2 **1 CHMP qualification Opinion statement ..... 3**

3 **2 Executive summary as submitted by the applicant..... 3**

4 **2.1 The objective(s) of the request.....3**

5 **2.2 The need and impact of proposed clinical novel methodology(ies) .....4**

6 **2.3 Sources of data .....7**

7 **2.4 Characteristics of the proposed novel methodology.....8**

8 **2.5 Differences between proposed COU and the Loupy et al., 2019 publication .....8**

9 **2.6 Summary of the Qualification purpose, methods, and results .....10**

10 **2.7 Overall goal of the present submission .....11**

11 **3 Questions from the Applicant and CHMP answers ..... 12**

12 **4 Background as submitted by the applicant..... 20**

13

14 Annexes to this Qualification Opinion published as separate documents as provided by the applicant:

15 • Validated Briefing Document providing background information

16 • Appendix to the Briefing Document (BD)

17 • Written Answers to List of Issues No. 1

18 • Written Answers to List of Issues No. 2

## 19 **1 CHMP qualification Opinion statement**

20 CHMP qualifies the iBox Scoring System (Composite Biomarker Panel) as a secondary endpoint  
21 prognostic for death-censored allograft loss (allograft failure) in kidney transplant recipients to be used  
22 in clinical trials to support the evaluation of novel immunosuppressive therapy applications.

23 This opinion applies to both the abbreviated and the full iBox Scoring System. Considering the  
24 minimal difference in the performance of these two scores and the requirement for an invasive  
25 procedure for the full iBox Scoring system, the abbreviated iBox Scoring System may be the preferred  
26 one. The scoring systems predict death censored allograft failure at 5 years. This is not the preferred  
27 primary clinical endpoint as the preferred primary estimand includes death as an observed event. This  
28 should be taken into consideration for the development of a surrogate endpoint and further work on an  
29 all-cause endpoint score should be pursued. It is acknowledged that prediction of all-cause death  
30 events may be challenging at an early time point post transplantation. Finally, in order to increase the  
31 number of trials fulfilling the criteria for validation studies, the Applicant should consider an outcome  
32 reflecting the assessment of efficacy of chronic kidney disease, i.e. relative reduction in eGFR (30 to  
33 57%) in addition to graft failure and death (EMA CKD guideline).

34 The focus of the analysis presented is to support use of the iBox score at 1 year post transplantation  
35 to assess 5-year risk of kidney graft failure. Nevertheless, the dataset supports a more flexible COU  
36 with the iBox score measured between 6- and 24-months post-kidney transplantation in pivotal or  
37 exploratory drug therapeutic studies for regulatory purposes. Additional material is provided to support  
38 this conclusion (Appendix to Briefing Document). The CHMP encourages the use of the iBox scoring  
39 system as a secondary endpoint in future trials of kidney transplantation and further development of  
40 the scoring system targeting a potential future qualification as a surrogate endpoint. Sponsors may  
41 consider using the iBox Scoring System as a secondary endpoint with Type 1 error control included in  
42 a procedure to address multiplicity in pivotal trials.

43 For a more detailed discussion of the CHMP assessment, please see '3. Questions posed by the  
44 applicant and CHMP answers'.

## 45 46 **2 Executive summary as submitted by the applicant**

### 47 **2.1 The objective(s) of the request**

48 The objective of this Briefing Dossier is for the Critical Path Institute's (C-Path) Transplant  
49 Therapeutics Consortium (TTC) to achieve a Qualification Opinion for a new drug development tool  
50 (DDT) for kidney transplantation through the EMA's qualification of novel methodologies for medicine  
51 drug development. This Briefing Dossier contains the proposed context-of-use (COU) statement, data  
52 source description, modeling analysis methods, and results that provide a quantitative basis to support  
53 the use of the iBox Scoring System (Composite Biomarker Panel), known as iBox Scoring System  
54 henceforth, as a surrogate endpoint for the five-year risk of death-censored allograft loss (allograft  
55 failure) in kidney transplant recipients for use in clinical trials evaluating the safety and efficacy of  
56 novel immunosuppressive therapies (ISTs). Two iBox Scoring System models have been developed  
57 and are included in this qualification submission: a full iBox Scoring System (with biopsy) and an  
58 abbreviated iBox Scoring System (without biopsy) known henceforth as the full iBox Scoring System,  
59 or the abbreviated iBox Scoring System, respectively. Additionally, a scoring system for predicting a  
60 combined endpoint including allograft failure and patient death as events), the ACE (all-cause  
61 endpoint) score, has been derived and tested in the external validation datasets

62 The iBox Scoring System has been developed by estimating individual weights for each of the  
63 proposed components (i.e., estimated glomerular filtration rate [eGFR] calculated by the 4-variable  
64 Modification of Diet in Renal Disease (MDRD)-186 Study equation, proteinuria, kidney allograft biopsy  
65 histopathology, presence of donor-specific antibodies [DSA], and time of post-transplant iBox Scoring  
66 System risk evaluation. For the purpose of this submission, the time of post-transplant risk evaluation  
67 was fixed at one-year post-transplant. The ACE score incorporates all of the variables in the  
68 abbreviated iBox Scoring System.

## 69 **2.2 The need and impact of proposed clinical novel methodology(ies)**

70 The two major transplantation societies in the United States, the American Society of Transplant  
71 Surgeons (ASTS) and the American Society of Transplantation (AST), recognized in 2014 the need for  
72 a pathway for the development of new ISTs for transplant recipients. (Stegall et al. 2016). The two  
73 societies, along with other members of the transplant community and C-Path, created the TTC. The  
74 goal of the TTC is the goal of this proposal—to develop a path forward to accelerate the medical  
75 product development process for transplantation, with a focus on novel ISTs that are likely to improve  
76 long-term renal allograft survival. Following the Loupy et al., 2019 publication introducing the iBox risk  
77 prediction tool, AST and ASTS signed a joint letter of support in March of 2020 encouraging the  
78 Institut national de la santé et de la recherche médicale (Inserm) to share patient-level data used to  
79 derive the iBox Scoring System as per Loupy et al., 2019 with the TTC. This letter of support was  
80 written to assist the regulatory endorsement of the iBox Scoring System as a surrogate endpoint in  
81 kidney transplant clinical trials. The joint letter of support can be found in Appendix (AST-ASTS TTC  
82 Joint Letter of Support).

83 The historically-accepted clinical trial endpoint for multinational clinical trials of novel ISTs in kidney  
84 transplantation is the composite endpoint of equally-weighted death, graft-loss, biopsy-proven acute  
85 rejection (BPAR) and lost to follow-up at one-year post-transplantation. There are several issues with  
86 the continued reliance on this endpoint with the current standard of care (SOC) ISTs. Firstly, the  
87 incidence is low in the first year post-transplant, limiting the ability to demonstrate the superiority of a  
88 new innovative therapy. Secondly, this endpoint was originally designed to quantify the incidence of  
89 BPAR without censoring. However, this approach results in the equal weighting of transplant recipients  
90 who die compared to those with BPAR or are lost to follow-up. Lastly, the largest unmet need in  
91 transplant is improvement in the long-term survival of the transplant recipient and graft and the  
92 associated surrogate endpoints that are predictive of survival. Current IST regimens have dramatically  
93 improved short-term outcomes, with one-year graft survival rates of approximately 91% after  
94 deceased donor transplant, according to the European Renal Association - European Dialysis and  
95 Transplant Association (ERA-EDTA) 2018 Annual Report (ERA-EDTA Registry Annual Report 2018).  
96 Despite these improved short-term outcomes, long-term graft survival remains suboptimal. The 5 -  
97 and 10-year graft survival rate after deceased donor kidney transplant is 77% and 56%, respectively  
98 (Gondos et al. 2013). Consequently, there is a significant unmet need for ISTs that can help improve  
99 long-term outcomes, but developing novel therapies is challenging. One aspect of this challenge is  
100 demonstrating improved long-term outcomes, which require trials of long duration (i.e., five years or  
101 more) and contain a large number of subjects. As a result, one-to-two-year non-inferiority studies are  
102 more likely to be initiated, despite not adequately addressing the challenges of improving long-term  
103 graft survival. A strategy of using surrogate endpoints in assessing long-term outcomes has been  
104 employed in other therapeutic areas, such as oncology, diabetes, nephrology, and many rare diseases,  
105 to overcome similar challenges. Surrogate endpoints enable sponsors to seek conditional marketing  
106 authorisation (CMA) for novel agents based on clinical trials of reasonable duration (i.e., one year) that  
107 predict long-term outcomes (i.e., five years or greater) while planning and conducting studies to  
108 demonstrate longer-term therapeutic effects.

109 The challenges associated with developing a robust surrogate endpoint capable of accurately predicting  
110 long-term outcomes (i.e., five-year risk of graft loss) using short-term data (i.e., one-year post-  
111 transplant) are multifaceted. Two of the most significant challenges include the need to develop a  
112 reliable surrogate measure that performs across a heterogeneous subject population and the ability of  
113 the surrogate measure to demonstrate efficacy across therapies with multiple mechanisms of action  
114 (MOA). In addition, subject-level data from various sources representing a broad spectrum of subject  
115 populations and treatment settings must be aligned and curated to generate the necessary evidence to  
116 support the surrogacy claims of such a measure.

117 In 2019, the Paris Transplant Group (French National Institute of Health), together with 29 key opinion  
118 leaders of the transplant community from 10 referral centers from Europe and the USA, published a  
119 seminal paper on the iBox Scoring System titled: Prediction system for risk of allograft survival in  
120 subjects receiving kidney transplants: international derivation and validation study (Alexandre Loupy  
121 et al. 2019). The PTG designed a prospective study to identify key prognostic parameters and follow  
122 long-term outcomes of kidney transplant recipients to develop a new risk prediction model of long-  
123 term kidney allograft failure outperforming previous scoring systems.

124 In this publication, the iBox Scoring System is a risk prediction tool utilizing multiple clinically relevant  
125 subject features of kidney function (eGFR and proteinuria), kidney allograft biopsy histopathology, and  
126 immunological status (presence of DSA) data cross-sectionally at any timepoint post-transplantation.  
127 The component measures of the iBox Scoring System are routinely used as important factors in  
128 routine monitoring of transplant recipients to guide therapeutic interventions and for prognostic  
129 purposes. The iBox Scoring System integrates these measures to generate individualized predictions of  
130 outcomes at three, five, and seven-years post-transplant. Data prospectively collected from 4,000  
131 consecutive subjects across four health centers in France were used to develop the iBox Scoring  
132 System, with external validation performed in cohorts from transplant centers in the U.S. (n = 1,428),  
133 Europe (n = 2,129), a phase III IST minimization trial (n = 194), a phase III trial assessing treatment  
134 of active antibody-mediated rejection (aAMR) in subjects with pre-transplant DSA (n = 38), and a  
135 phase II trial evaluating treatment of antibody-mediated rejection (AMR) in subjects with post-  
136 transplant de novo DSA (n = 44). The TTC, in close collaboration with the PTG, is seeking to translate  
137 the work from Loupy et al., 2019 British Medical Journal (BMJ) publication into a regulatory endpoint in  
138 hopes of streamlining drug development by facilitating clinical trials of shorter duration (i.e., one year)  
139 that can predict death-censored allograft survival.

140 While the underlying physiological mechanisms leading to allograft survival are complex, recent  
141 studies have shown that certain key features present relatively early after transplantation (i.e., within  
142 the first year) can accurately predict which grafts are most likely to fail at later time points (i.e., by  
143 five years). A key learning from prior efforts in the field is no one clinical feature or pathophysiological  
144 measure has the predictive power to robustly estimate long-term allograft survival (Naesens et al.  
145 2016); (Kaplan, Schold, and Meier-Kriesche 2003); (Yilmaz et al. 2003); (Lefaucheur et al. 2010).  
146 Recent efforts that have had access to large subject cohorts with rigorous and routine clinical  
147 assessments collected at baseline and longitudinally for five to seven years have demonstrated  
148 improved predictability of long-term outcomes by assessing composites of multiple clinical features.  
149 These composite scores have focused on recipient demographics, pre-transplant measures, measures  
150 of kidney function within the first-year post-transplant, and combinations of these measures at  
151 different time points (Kaboré et al. 2017); (Shabir et al. 2014); (Gonzales et al. 2016); (Alexandre  
152 Loupy et al. 2019);(Rampersad et al. 2021).

153 More recently-developed composite scores have sought to predict long-term graft loss by incorporating  
154 a cross-section of the relevant pathophysiological measures of allograft survival, including kidney  
155 function, through eGFR calculated using serum creatinine (SCr) and measures of protein excreted into

156 the urine, kidney damage as determined by pathological assessment of graft biopsy, and immune  
157 response, measured via the presence of DSA. Other composite scores have incorporated  
158 pathophysiological measures and recipient demographics (Gonzales et al. 2016); (Bentall et al. 2019).

159 These risk prediction scores have focused on predicting long-term allograft survival at the subject-level  
160 to inform individual clinical decision-making. However, none of these tools have been subject to  
161 independent external validation. Consequently, none of these tools have been a candidate or endorsed  
162 for use as a surrogate endpoint capable of supporting medical product registration studies or as  
163 surrogate endpoints in the context of EMA's CMA (Menon, Murphy, and Heeger 2017). On the contrary,  
164 the iBox Scoring System showed accuracy in predicting death-censored allograft failure, which was  
165 confirmed across transplant centers worldwide, different subpopulations and clinical scenarios, as well  
166 as in randomized clinical trials (RCTs), lending its exportability to a variety of clinical trial settings.

167 The proposed iBox Scoring System in this submission is intended to be a surrogate endpoint for  
168 efficacy in clinical trials evaluating the safety and efficacy of novel ISTs in kidney transplant recipients  
169 as a marker for the probability of long-term allograft survival. TTC aims to improve upon the  
170 limitations of the historically utilized clinical trial primary endpoint through the development and  
171 regulatory endorsement of the iBox Scoring System capable of predicting long-term kidney transplant  
172 outcomes using measures available at one-year post-transplantation.

173 This effort builds on previous work in the field that has identified clinically relevant measures capable  
174 of predicting long-term allograft failure by curating data from multiple clinical trials, real-world clinical  
175 transplant center datasets, and long-term registry data. The TTC has been working closely with the  
176 PTG and the global transplant community to curate and align subject-level data to support the use of  
177 the iBox Scoring System in drug development. A key difference between the iBox Scoring System in  
178 the Loupy et al., 2019 manuscript and the iBox Scoring System as a surrogate endpoint detailed in  
179 this submission, is the time point for risk evaluation. In this submission, the COU has been defined  
180 with the risk evaluation fixed at one year post kidney transplant. While the Loupy, et al., 2019 iBox  
181 Scoring System algorithm allows the risk to be estimated at any time point post-transplant. The COU  
182 in this submission prespecified the risk evaluation at one-year post-transplant to adapt the iBox  
183 Scoring System described in Loupy et al. into a clinical trial endpoint at a fixed landmark. In order to  
184 facilitate the use of the iBox Scoring System in a multinational clinical trial, two versions of the iBox  
185 Scoring System were assessed, one version including all components as described by Loupy et al.,  
186 2019 (Full iBox Scoring System) and one version excluding pathophysiological assessment of the  
187 kidney allograft biopsy (abbreviated iBox Scoring System). Also, to adapt the Loupy et al., 2019 iBox  
188 Scoring System to be used as a one-year clinical trial endpoint, analyses were performed imputing a  
189 one-year iBox score for subjects who died or lost a graft in the first-year post-transplant.

190 Based on existing literature and work by the PTG, the proposed components of the iBox Scoring  
191 System model include:

- 192 • eGFR calculated by the 4-variable MDRD-186 Study equation with SCr (referred to as 'eGFR');
- 193 • Measurement of protein excretion into the urine through calculation of the urine protein-to-  
194 creatinine ratio (referred to as 'proteinuria');
- 195 • Histopathological assessment of tissue obtained by renal allograft biopsy (referred to as  
196 'kidney allograft biopsy histopathology');
- 197 • Presence of DSA;
- 198 • The time of post-transplant iBox Scoring System risk evaluation. For the purpose of this  
199 submission, the time of risk evaluation was fixed at one-year post-transplant.

200 The multivariable Cox PH model was used to adapt the full and abbreviated iBox Scoring System  
201 models for use at one-year post-transplant as a surrogate endpoint for the five-year risk of death-  
202 censored allograft survival. Thus, this Briefing Dossier will consist of a discussion of these proposed  
203 components.

### 204 **2.3 Sources of data**

205 To acquire the subject-level data necessary to develop a novel surrogate endpoint, the TTC led an  
206 extensive global data collaboration effort across the field of kidney transplantation. To date, the TTC  
207 has acquired eleven clinical trial datasets and twenty observational datasets from clinical transplant  
208 centers, representing data from over 20,000 kidney transplant recipients in the TTC Kidney Transplant  
209 Database. A list of acquired datasets can be found in the Appendix (Revised-Transplant Therapeutics  
210 Consortium's Kidney Transplant Database).

211 Datasets from relevant clinical trials of ISTs, including those in the Loupy et al. 2019 publication, and  
212 real-world data from international clinical transplant centers were prioritized for acquisition. From  
213 these 31 datasets, five contained all necessary variables collected at one-year post-transplant (i.e.,  
214 eGFR, proteinuria, kidney allograft biopsy histopathology, and DSA), long-term death and graft loss  
215 follow-up of at least five years, immunosuppressive regimen information (i.e., induction and  
216 maintenance IST) to test the performance of the surrogate with all three MOA, and the documentation  
217 required to support the description of the analytical considerations for each dataset.

218 Datasets missing the necessary variables at one-year post-transplant or a variable necessary to  
219 calculate the model variable (as in recipient age to calculate an eGFR value) were excluded. For  
220 example, in the data for the three Novartis studies (TRANSFORM, US-92, and ELEVATE), recipient age  
221 was missing due to Novartis' anonymization procedures for data sharing. This, in turn, prohibited the  
222 calculation of eGFR values for the subjects in these studies. Moreover, US-92 and ELEVATE were  
223 missing DSA and proteinuria data, and follow-up was limited to one and two years, respectively.

224 The five datasets described below were therefore used for this qualification submission. These five  
225 qualification datasets consist of one derivation dataset and four validation datasets, outlined below.

#### 226 **Qualification derivation dataset:**

227 1. The qualification derivation dataset presented in this Briefing Dossier included specific  
228 adjustments to the original derivation dataset as described in Loupy et al., 2019 manuscript,  
229 (Alexandre Loupy et al. 2019), allowing the iBox Scoring System to be used as a one-year  
230 post-transplant surrogate endpoint in clinical trials. This data was received from the PTG in  
231 Paris, France, Europe consisting of the following four transplant centers:

- 232 • Necker Hospital in Paris, France, Europe.
- 233 • Saint-Louis Hospital in Paris, France, Europe.
- 234 • Foch Hospital in Suresnes, France, Europe.
- 235 • Toulouse Hospital in Toulouse, France, Europe.

#### 236 **Qualification validation datasets:**

237 The qualification validation datasets presented in this Briefing Dossier contain datasets other than  
238 those used for external validation as described in Loupy et al., 2019 manuscript (Alexandre Loupy et  
239 al. 2019). The qualification validation datasets are from both transplant centers and RCTs as described  
240 below.

241

- 242 2. Mayo Clinic in Rochester, Minnesota, USA, North America.
- 243 3. Helsinki University Hospital in Helsinki, Finland, Europe.
- 244 4. A phase III study of belatacept-based immunosuppression regimens versus cyclosporine (CsA)
- 245 in recipients of kidneys from living or standard criteria deceased donor kidneys (BENEFIT RCT)
- 246 Vincenti et al., 2012.
- 247 5. A phase III study of belatacept versus CsA in recipients of allografts from extended criteria
- 248 donors, those donated after cardiac death, and those with an estimated cold ischemia time
- 249 (CIT) > 24 hours in duration (BENEFIT-EXT RCT) Medina-Pestana., 2012

250 The qualification derivation and validation datasets were aligned and curated to support the regulatory

251 endorsement of the full and abbreviated iBox Scoring System models. These datasets were used to

252 construct the statistical analysis plan (SAP) presented in this Briefing Dossier.

## 253 **2.4 Characteristics of the proposed novel methodology**

### 254 **Proposed context-of-use statement**

255 The iBox Scoring System (Composite Biomarker Panel) used at one-year post-transplant is a surrogate

256 endpoint for the five-year risk of death-censored allograft loss (allograft failure) in kidney transplant

257 recipients for use in clinical trials to support evaluation of novel IST applications via CMA.

#### 258 **General area:**

259 Surrogate endpoint for the five-year risk of death-censored allograft loss (allograft failure) in kidney

260 transplant subjects for use in clinical trials to support evaluation of novel IST applications.

#### 261 **Target population for use of the biomarker:**

262 Adult *de novo* kidney only transplant recipients from a living or deceased donor.

#### 263 **Stage of drug development for use:**

264 All clinical efficacy evaluation stages of therapeutic interventions focused on the use of the long-term

265 risk of allograft survival in kidney transplant recipients, including early signs of efficacy, proof-of-

266 concept, dose-ranging, and registration studies (Phases II-IV).

#### 267 **Intended application:**

268 The iBox Scoring System (Composite Biomarker Panel) used at one-year post-transplant is a surrogate

269 endpoint for the five-year risk of death-censored allograft loss (allograft failure) in kidney transplant

270 subjects for use in clinical trials to support evaluation of novel IST applications via CMA. When

271 evaluating five-year outcomes for clinical benefit and full marketing authorisation, it will be necessary

272 to ensure that there is not a clinically meaningful decrease in transplant recipient survival with the new

273 therapy in the clinical trial compared to the standard of care control arms.

## 274 **2.5 Differences between proposed COU and the Loupy et al., 2019**

### 275 **publication**

276 The original derivation dataset (Alexandre Loupy et al. 2019) was used in the derivation analysis of the

277 full iBox Scoring System and the abbreviated iBox Scoring System. The qualification derivation dataset

278 presented in this Briefing Dossier included specific adjustments to the originally derived formula

279 allowing the iBox Scoring System risk evaluation at one-year post-transplantation for use in a clinical

280 trial endpoint at a fixed landmark. The qualification validation presented in this Briefing Dossier used



281 datasets other than those used for external validation in Loupy et al., 2019 manuscript [(Alexandre  
282 Loupy et al. 2019).

283 Table 1. compares and contrasts the iBox Scoring System described in Loupy et al., 2019 manuscript  
284 and the iBox Scoring System as a surrogate endpoint proposed in this Briefing Dossier for Qualification  
285 Opinion.

286 **Table 1. iBox Scoring System as described in Loupy et al., 2019 versus iBox Scoring System**  
287 **proposed for Qualification Opinion**

	Loupy et al., 2019	iBox Scoring System proposed for Qualification Opinion
<b>Core components of model</b>	<ol style="list-style-type: none"> <li>1. eGFR<sub>MDRD</sub></li> <li>2. Proteinuria: log transformed UPCR</li> <li>3. Kidney allograft biopsy histopathology</li> <li>4. DSA: Semiquantitative mean fluorescence intensity (MFI) associated with DSA</li> <li>5. Time of post-transplant risk evaluation: at any time from transplant</li> </ol>	<ol style="list-style-type: none"> <li>1. eGFR<sub>MDRD</sub></li> <li>2. Proteinuria: log transformed UPCR; imputation methodology included for datasets using other proteinuria measurements</li> <li>3. Two iBox Scoring System models, one with and one without kidney allograft biopsy histopathology</li> <li>4. DSA: Binary qualitative MFI associated with DSA</li> <li>5. Time of post-transplant risk evaluation: one-year post-transplant</li> </ol>
<b>Application</b>	Individual decision-making	Surrogate endpoint in kidney transplantation clinical trials
<b>Derivation set</b>	Loupy et al., 2019	Loupy et al., 2019
<b>External validation sets</b>	Hôpital Hôtel Dieu, Nantes, France; Hospices Civils, Lyon, France; University Hospitals, Leuven, Belgium; Johns Hopkins Medical Institute, Baltimore, MD; the Mayo Clinic, Rochester, MN; and the Virginia Commonwealth University School of Medicine, Richmond, VA	Mayo Clinic Rochester <sup>†</sup> ;  Helsinki University Hospital;  BENEFIT RCT;  BENEFIT-EXT RCT
<b>Methodology</b>	Semiparametric Cox PH model	Semiparametric Cox PH model; imputation for proteinuria and for subjects who die or lose their graft in the first year of transplant

<b>Outcomes</b>	Death-censored allograft survival	Death-censored allograft survival
<b>Imputation used for sensitivity analysis in trial-level surrogacy (TLS) and for one-year endpoint definition</b>	No	Yes
<b>Assay documentation</b>	Excluded	Included

288 † Different dataset than in Loupy et al., 2019

## 289 **2.6 Summary of the Qualification purpose, methods, and results**

290 There is a need for new short-term endpoints in kidney transplant trials that allow demonstration of  
 291 superiority of new therapies over the current SOC and translate into reductions in long-term graft loss.  
 292 The availability of a surrogate endpoint is vital to stimulate innovation in immunosuppressive drug  
 293 development that will serve transplant recipients by further improving short- and long-term outcomes.

294 Loupy et al., 2019 developed the iBox Scoring System as a risk prediction score for death-censored  
 295 kidney allograft survival by estimating individual weights for each of the proposed components (i.e.,  
 296 eGFR, proteinuria, kidney allograft biopsy histopathology, the presence of DSA, and time of post-  
 297 transplant risk evaluation). The TTC has adapted the innovative work by Loupy et al., 2019, to  
 298 transform the original iBox Scoring System to a surrogate clinical trial endpoint measured at one-year  
 299 post-transplant.

300 The following key analyses have been performed and are detailed in this submission:

- 301 • Original iBox Scoring System analyses of data by Loupy et al., 2019 have been reproduced for  
 302 the full iBox Scoring System and abbreviated iBox Scoring System for the data from the PTG  
 303 (derivation dataset n = 3,941 for full iBox Scoring System and n = 4,000 for abbreviated iBox  
 304 Scoring System).
- 305 • For application as an endpoint in a clinical trial at one-year, the derivation dataset from PTG  
 306 was analyzed, restricting the analysis to those recipients with a full iBox Scoring System  
 307 evaluation at one-year post-transplant and follow-up to five-years for graft loss (n = 1,174).  
 308 The discrimination in this group was confirmed with a c-statistic = 0.85.
- 309 • Subsequently, external validation was performed in the four qualification datasets (i.e., two  
 310 observational datasets from Helsinki University Hospital and Mayo Clinic Rochester and two  
 311 RCTs from Bristol-Meyers Squibb (BMS), BENEFIT and BENEFIT-EXT).
- 312 • External validation was performed using discrimination (c-statistics) and calibration (observed  
 313 versus predicted graft loss) methods. In all four of the qualification validation datasets using  
 314 the full and abbreviated iBox Scoring System models at one year to predict five-year death-  
 315 censored allograft survival, the c-statistics ranged from 0.70-0.93, and the predicted versus  
 316 observed graft losses were not significantly different. These data confirmed the external  
 317 validation of the full and abbreviated iBox Scoring System.

- 318 • Discrimination (c-statistics) was also included for the European validation cohort (c-statistic =  
319 0.81, 95% confidence interval [CI] 0.78 to 0.84) and the three RCTs, [CERTITEM (c-statistic =  
320 0.88), RITUX ERAH (c-statistic = 0.77), and BORTEJECT (c-statistic = 0.94)] described in  
321 Loupy et al., 2019 as additional data supporting this qualification submission.
- 322 • The ability of the iBox Scoring System to demonstrate a treatment effect at one-year that  
323 translates into a treatment effect on death-censored five-year graft survival was assessed in  
324 two ways. First, TLS was performed but, due to insufficient data (i.e., only two prospective  
325 RCTs and a mTORi derivation subset), it was not possible to provide the precise estimation of  
326 the trial-level correlation coefficient. Study level treatment effects in the BENEFIT RCT,  
327 BENEFIT EXT RCT, and a mTORi derivation subset using a calcineurin inhibitor (CNI) free  
328 regimen, mammalian target of rapamycin (mTORi) such as sirolimus or everolimus versus  
329 CNI-based regimen data from Loupy et al., 2019 qualification derivation data for one-year iBox  
330 scores for the full and abbreviated iBox Scoring System and five-year death-censored allograft  
331 survival were also assessed. The average iBox score at one year was consistently significantly  
332 lower in the CNI-free arm (belatacept [BELA] or mTORi) compared to CNI arms. The five-year  
333 death-censored allograft survival also consistently numerically favored the CNI-free arm. At  
334 five-years in the BENEFIT RCT, death-censored allograft survival was significantly better with  
335 BELA compared to CsA. Analyses of the BENEFIT RCT included imputation of the worst-case  
336 iBox Scoring System at one-year post-transplant for recipients who died or lost their graft in  
337 the first year. This sensitivity analysis was performed to replicate the clinical trial setting  
338 where avoidance of survivor bias at one year would be necessary, and all randomized subjects  
339 would have an iBox score at one-year even if there were death or graft loss before that time.  
340 The totality of these data demonstrate that the iBox Scoring System can measure treatment  
341 effects at one-year that translate into a consistent impact on the five-year death-censored  
342 allograft survival. The lack of statistical significance on some of the five-year death-censored  
343 allograft survival analysis is related to limitations in power to detect differences based on  
344 sample size.

345 Based on these analyses, the full or abbreviated iBox Scoring System models at one-year post-  
346 transplant is a validated surrogate for the five-year death-censored allograft survival and is applicable  
347 for use in a prospective RCT with imputation for deaths and graft losses within the first year of  
348 transplant. Qualification of the iBox Scoring System as a surrogate endpoint would significantly  
349 improve upon the current standard, as it would allow drug sponsors the ability to design trials  
350 assessing the superiority, of a novel agent. As a surrogate endpoint for the long-term outcome of  
351 allograft survival, the iBox Scoring System would allow drug sponsors to seek marketing authorisation  
352 of novel agents through EMA's CMA process while planning and conducting additional studies to  
353 demonstrate longer-term therapeutic effects, thereby significantly improving the drug development  
354 landscape by encouraging drug sponsors to engage in this therapeutic area of high unmet need.  
355 Ultimately, kidney transplant recipients will benefit from the increased drug development activity by  
356 improving access to ISTs with better short-term and long-term outcomes.

## 357 **2.7 Overall goal of the present submission**

358 The TTC presents this Briefing Dossier to request a Qualification Opinion from the Agency on the  
359 proposed COU for the iBox Scoring System at one-year post-transplant as a surrogate endpoint for the  
360 five-year risk of death-censored allograft loss (allograft failure) in kidney transplant subjects for use in  
361 clinical trials to support evaluation of novel IST applications via CMA process. The TTC believes a  
362 Qualification Opinion is critical for accelerating the development of ISTs in kidney transplantation  
363 clinical trials.

### 364 **3 Questions from the Applicant and CHMP answers**

#### 365 **Does EMA agree with the COU?**

366 **TTC's position:** The proposed COU provides a quantitative basis to support the use of the iBox  
367 Scoring System (Composite Biomarker Panel) at one-year post-transplant as a surrogate endpoint for  
368 the five-year risk of death-censored allograft loss (allograft failure) in kidney transplant recipients for  
369 use in a clinical trial endpoint at a fixed landmark. Qualifying two iBox Scoring System models, with  
370 and without biopsy input, will provide sponsors and investigators flexibility in clinical trial design, with  
371 or without a surveillance biopsy at one-year post-transplant.

372 As this surrogate endpoint is proposed to be used in the context of CMA with EMA, where full approval  
373 of a product will not be authorized until the clinically meaningful outcome (five-year death-censored  
374 allograft survival) has been met, the TTC feels that sufficient evidence is provided in this dossier to  
375 support qualification of the iBox Scoring System.

#### 376 **CHMP answer**

377 It is agreed that there is a need to develop a reliable surrogate measure that performs across a  
378 heterogenous population and allow showing efficacy across therapies with multiple mechanisms of  
379 action (MoA).

380 The initially proposed Context of Use (COU) for the two composite biomarker panels was use as a  
381 surrogate endpoint to predict the five-year risk of death-censored allograft loss (allograft failure) in  
382 kidney transplant recipients for use in clinical trials to support evaluation of novel IST applications via  
383 CMA. The target population are adult de novo kidney only transplant recipients from a living or  
384 deceased donor. Development of two scores (one with and one without histology) seems reasonable  
385 and could ease recruitment and maintenance of patients in future studies; biopsy may be associated  
386 with bleeding, renal fistulas and haematuria. The transplant recipient may refuse biopsy for study  
387 purposes only.

388 After discussion of two lists of issues provided by SAWP, the COU was modified and refined with a final  
389 proposal of the statement reading 'The iBox Scoring System (Composite Biomarker Panel) is a co-  
390 primary or secondary endpoint prognostic for death-censored allograft loss (allograft failure) in kidney  
391 transplant recipients to be used in clinical trials to support the evaluation of novel immunosuppressive  
392 therapy applications.' Additional information was provided by the Applicant that supports a more  
393 flexible COU with the iBox measured between 6- and 24-months post-kidney transplantation in pivotal  
394 or exploratory drug therapeutic studies for regulatory purposes. While the focus would likely be long-  
395 term prediction of death-censored graft loss, also shorter periods for prediction would be feasible with  
396 less events expected in a shorter time frame. From regulatory point of view the preferred primary  
397 clinical endpoint is to include death as an observed event. This should be taken into consideration for  
398 future development of the iBox Scoring System.

399 The more flexible COU would allow using the iBox scoring system in proof of concept or dose finding  
400 phase 2 studies and phase 3 studies. It is possible that iBox could add supportive evidence for CMA,  
401 provided requirements for CMA are fulfilled. These are outlined in EMA guideline  
402 (EMA/CHMP/509951/2006, Rev.1). For the iBox to support CMA, it will have to be able to support a  
403 positive benefit-risk balance of the medicine under investigation and it will have to be ensured that it  
404 is likely that comprehensive data post-authorisation will be generated. The timeframe to provide data  
405 post-authorization should not jeopardize the conduct of the study, e.g., in case of availability of a  
406 newly approved medical therapy.

407 The Applicant states (chapter 3.2) that when using the death-censored iBox score, it will always be  
408 necessary to determine if there is clinically meaningful decrease in transplant recipient survival with  
409 new therapy. This view is shared. Other post-authorisation requirements for CMA include the fulfilment  
410 of an unmet medical need and the benefit of the medicine's immediate availability to patients is  
411 greater than the risk inherent in the fact that additional data are still required.

412 There are several other regulatory approaches available to address safety, and/or efficacy, post  
413 approval. Such, post-authorisation measures (PAMs) may be aimed at collecting or providing data to  
414 enable the assessment of the safety or efficacy (see EMA website "[Post-authorisation measures:  
415 questions and answers](#)").

416 In conclusion, the initial COU proposed for iBox scoring system was a surrogate endpoint to support  
417 CMA. As explained, a surrogate endpoint is not a priori linked to a specific regulatory pathway within  
418 the EU. Please see the discussion regarding the assessment of iBox scoring system as a surrogate  
419 endpoint in the answer to Q3.

420

#### 421 **Does EMA agree that the data sources are adequate to support the proposed COU?**

422 **TTC's position:** The TTC led an extensive data collaboration effort across the field of kidney  
423 transplantation. Datasets from relevant clinical trials of ISTs, including the data in Loupy et al., 2019  
424 publication and real-world data from international clinical transplant centres, were prioritized. There  
425 were five datasets that contained all of the necessary clinical variables collected at one-year post-  
426 transplant (i.e., eGFR, proteinuria, kidney allograft biopsy histopathology, and presence of DSA), long-  
427 term death and graft loss follow-up of at least five years, immunosuppressive regimen information  
428 (i.e., induction and maintenance IST) to test the performance of the surrogate with all three MOA, and  
429 the documentation required to support the description of the analytical considerations for each dataset  
430 in this qualification submission. C-Path has reviewed the documentation and deemed that the  
431 analytical methods were robust, reliable, and fit-for-purpose.

432 The available data sources, and their alignment through experienced and quality data management,  
433 represent a unique opportunity to transform these data into valuable knowledge to provide the  
434 necessary evidence to support the qualification of the iBox Scoring System (Composite Biomarker  
435 Panel) for the proposed COU. The population captured in the data sources represents the population  
436 likely to be considered as candidates to participate in clinical trials of therapies intended to improve  
437 long-term graft survival.

#### 438 **CHMP answer**

439 It is agreed that the clinical transplant population is heterogenous. This also poses a challenge to  
440 establishing surrogacy. The proposed target population is "Adult *de novo* kidney only transplant  
441 recipients from a living or deceased donor", i.e. the broad population of adult transplants. The efforts  
442 of the TTC to acquire subject-level data for development of the proposed surrogate endpoint are  
443 acknowledged. Selecting studies (five out of 31) which contained all variables of interest is a  
444 reasonable approach. The variables with the composite panel are clinically relevant as they provide  
445 information on the health of the graft through measuring of renal function (proteinuria, eGFR), direct  
446 assessment of allograft health through histopathology, and the patient's immune response (DSA).

447 The five qualification datasets consist of one derivation dataset and four validation datasets; the latter  
448 comprised two prospective RCTs (the BENEFIT study and the BENEFIT EXT with a different target  
449 population). Subjects with grafts that never functioned (primary non-function) were excluded from the  
450 derivation data set. The broad range of patients and the variety data sources in the data set are

451 acknowledged. The documentation of the laboratory assays used is adequate and supports reliability  
452 and adequacy of the analytical laboratory methods. There was reclassification applied to address the  
453 fact that different criteria for graft loss were used across the data sets. This led to a number of  
454 reclassifications and there was a considerable number in the BENEFIT and BENEFIT-EXT studies.  
455 Standardisation of criteria, using consensus criteria according to Levin et al. (Levin A et al., Kidney  
456 International 2020) was implemented during the validation procedure, but is in principle welcomed and  
457 obviously important. The ad hoc reclassification was discussed at the first discussion meeting (DM) and  
458 there was no impact on interpretation of calibration results.

459 The studies included in the qualification exercise represent subjects with varying underlying diagnoses,  
460 receiving living related as well as extended donor kidneys, receiving various induction therapies and  
461 either CNI or CNI free therapy. As such, the data sources included are generally acceptable. However,  
462 the size of the database of the external validation studies is too small to determine consistency of the  
463 data across subpopulations. Also, most of these subsets are limited to single treatment centres. A  
464 limitation of the data sources is the small number of patients included in therapeutic intervention trials  
465 that are important for assessing the change in treatment effects in the proposed surrogate and the  
466 clinical endpoint at 5 years. Outcome events derived from randomised controlled trials are too sparse  
467 to be fully informative for the surrogacy at trial level of the iBox biomarker panel. To illustrate, in the  
468 largest trial there were 416 subjects with full iBox data at one year in the BENEFIT RCT, of whom 12  
469 graft losses at 5 years were recorded.

470 The low number of endpoint events in the available trials with an intervention limit establishing a  
471 correlation of treatment induced modification of the surrogate to treatment induced modification of the  
472 endpoint at 5 years. Such a relation is considered key for establishing full surrogacy of a biomarker-  
473 based endpoint. The correlation coefficients indicating the relation between treatment effect on iBox  
474 score and treatment effect on 5-year allograft survival were positive but low (0.0307-0.3054).  
475 Splitting the data into pseudo-trials per region as performed by the Applicant (p. 121 BD) was helpful  
476 in allowing further assessment of the correlation at (pseudo-)trial level but does expectedly not  
477 contribute much to improve precision of estimates for elucidating trial level surrogacy. Trial level  
478 surrogacy is assessed in the answer to Question 3. The ongoing efforts of the TTC to explore if  
479 additional RCTs exist that that may support the trial-level surrogacy (TLS) are acknowledged. The  
480 notion that there are insufficient completed RCTs in existence globally to execute a reasonable TLS  
481 analysis is noted.

482 During the DM several approaches were discussed which would potentially increase the number of  
483 trials fulfilling the criteria for validation of the iBox. These include using clinical trials that do not collect  
484 histology results to at least validate the abbreviated iBox, using outcome data at 3 years following  
485 transplantation, and redefining the outcome data to include relative reduction in eGFR. However, as  
486 per the Applicant, none of these measures were found to improve the number of trials available for  
487 validation of the iBox.

488 Taken together, the whole exercise would benefit from access to more data. Extensive global effort to  
489 collect clinical trials and real-world data on the side of the Applicant is understood and appreciated.

490

491

492 **Does EMA agree that the iBox Scoring System (Composite Biomarker Panel) or the all-cause**  
493 **endpoint (ACE) score have been validated as a surrogate endpoint for use in CMA**  
494 **submissions per their respective COU?**

495 **TTC's position:** The iBox Scoring System has been internally validated by the PTG and externally  
496 validated based on data from two transplant centres (one in Europe and one in the USA) and two  
497 Phase III multicentre, multinational RCTs. This external validation demonstrated both calibration and  
498 discrimination across the four qualification datasets. The presented analyses show that the iBox  
499 Scoring System can discriminate between higher and lower risk subjects in diverse datasets, including  
500 CNI and CNI-free populations. The results also showed the full and abbreviated iBox Scoring System  
501 had good prediction accuracy based on calibration analysis, including CNI and CNI-free populations in  
502 both transplant centres and RCTs.

503 The presented results demonstrate that the full and abbreviated iBox Scoring System models at one-  
504 year post-transplant are validated surrogates for the five-year death-censored graft survival and are  
505 applicable for use in a prospective RCT with imputation for deaths and graft losses within the first year  
506 of transplant.

507 The iBox Scoring System was designed to assess the long-term risk of allograft failure. Graft failure is  
508 defined as return to dialysis or pre-emptive re-transplantation. Death of the recipient with a  
509 functioning graft is typically a primary safety endpoint, with a wide variety of underlying causes of  
510 death observed (e.g., malignancy, infection, cardiovascular disease) and different risk factors  
511 compared with those for graft failure.

512 The ACE score has been internally validated in the qualification derivation dataset and externally  
513 validated in the qualification validation datasets. The ACE score was found to have modest  
514 discrimination, calibration, and predictive ability of a treatment effect in *de novo* kidney transplant  
515 recipients when high-risk donors were excluded and reduced discrimination as compared to the iBox  
516 Scoring System for predicting allograft loss.

#### 517 **CHMP answer**

##### 518 Overall validation approach

519 CHMP acknowledges the strengths of the current model development and validation approach and also  
520 the extensive and valuable initial work of the group led by Loupy (Loupy A et al, BMJ 2019). The initial  
521 prospective approach by Loupy et al. for derivation data collection led to a prediction model has good  
522 predictive performance for clinical endpoint events based on a number of variables included in a  
523 biomarker panel proposed as iBox. The model was internally and externally validated. Based on  
524 feedback from CHMP in a scientific advice on a proposal to use iBox as surrogate endpoint in a clinical  
525 phase 3 trial (EMA/CHMP/SAWP/650635/2019), the Applicant refined the approach and performed  
526 additional analysis. Inclusion of a new independent set of validation data for the refined development  
527 of the iBox score by TTC is welcomed by CHMP.

528 The differences to the initial approach by Loupy et al. are comprehensively explained in the BD (p.  
529 21). These include a different approach to handling donor specific antibodies (DSA) and pertain to the  
530 fixed 1-year time point proposed by the Applicant for the COU, which was addressed by imputing data  
531 for patients who die or lose graft during the first year. Imputation or spot proteinuria to reflect UPCR  
532 was performed for three of the four validation studies (BENEFIT, BENFIT-EXT and population from  
533 Helsinki University Hospital) and discussed below.

534 Two iBox models are proposed and this is in principle supported to allow flexibility in application in trial  
535 settings. The abbreviated iBox without biopsy information is supported by only a minimally larger  
536 number of subjects in the derivation data set (n=4000 vs. n=3941) and was retained after dropping  
537 the four kidney allograft biopsy histopathology variables in Table 38 of the BD. Backward elimination  
538 was not repeated after dropping the biopsy variables; the main goal with the abbreviated iBox was  
539 showing that dropping biopsy variables had minimal impact on model performance in the external

540 datasets. In the external validation dataset, more data without biopsy information are available. The  
541 development approach of the abbreviated model was discussed at the first meeting, e.g., if an  
542 abbreviated iBox could be re-derived with omitting biopsy related information. The Applicant explained  
543 that the 31 candidate variables explored in the derivation of the iBox Scoring System are not  
544 consistently present in the qualification validation datasets. It can also be concluded that restricting  
545 the analysis to an abbreviated iBox Scoring System will not increase the available data for analyses.

546 Missing data is minimal in derivation data set, any covariate imputation approach (opposed to  
547 imputation of iBox for patients that do not reach the 1-year time point) has no considerable input.  
548 Model development and analysis for internal validation was mainly data driven and this is acceptable in  
549 the given setting with pre-planned external validation based on additional independent data sets. The  
550 step of establishing trial level surrogacy for full validation of iBox has limitations, mainly due to the  
551 available datasets with a low number of observed events (please see below).

#### 552 Modelling and statistical methods

553 The modelling approach can be endorsed. As primary event of interest, graft loss was defined and  
554 death and loss to follow up were censored, assuming that these events are non-informative. As death  
555 as competing event could be informative, a competing risks analysis was performed. This is considered  
556 adequate. Subjects who died/withdrew/lost their graft before the first year after transplantation have  
557 missing iBox score values. These subjects were assigned imputed iBox score values. This is deemed a  
558 reasonable approach to avoid survivor bias. Incorporating scores for subjects who died for application  
559 of iBox with censoring for death using worst case scenario values for iBox at 1 year can be supported  
560 in principle (p. 56 BD).

561 A separate modelling approach using all-cause graft survival was also pursued to assess model  
562 performance when death is included in the model. The process of model derivation is appropriate.  
563 Univariate and multivariate analysis was used for variable selection from the 31 candidate variables.  
564 Backward elimination due to clinical considerations and rationale for categorical breakdown of  
565 variables in the univariate and multivariate models is comprehensively explained and can be  
566 supported. Overall, model assumption assessment for the Cox proportional hazard model and  
567 assessments of linearity of covariate using martingale residuals are endorsed. Testing the  
568 discriminatory properties for patients with and without graft loss (e.g., by ROC curve, p. 78 BD) is  
569 considered adequate. Using log transformed proteinuria values due to skewed data distribution  
570 appears adequate.

571 For performance assessment, Harrell's c-statistic was used (Harrell F, Stat Med 1996). This is an  
572 appropriate metric. Based on this measure, performance in patients without CNIs was assessed, as the  
573 training data used mainly subject treated with CNIs. Additionally, model performance with center as  
574 stratification factor was explored. Both steps are adequate and contribute to the validation. For an  
575 assessment of the predictive properties of the model with regard to accurately predicting the absolute  
576 risks, for calibration the number of predicted clinical endpoints were derived based on a Poisson model  
577 and compared to the observed events (Crowson C et al., Stat Methods Med Res 2016). The method of  
578 assessment of calibration is supported.

579 Supplementary assessment to assess the proteinuria conversion, death as competing risk for graft loss  
580 in the death-censored model and trial level surrogacy was performed. All these analyses are  
581 appropriately conducted and comprehensively described. It should be noted that imputation of urine-  
582 dipstick reflects spot concentration of urine-albumin and may change, e.g., with increased fluid intake,  
583 which is not the case for 24-hour proteinuria or UPCR. It is understood that the extrapolation of spot  
584 urine albumin by dipstick was based on a German population with both UPCR and dip stick results. The  
585 approach was further discussed at the first discussion meeting, as fit of the data was not clearly



586 demonstrated and the IQR (middle 50%) presented is very wide (figure 16). It can be concluded from  
587 the results and discussion that the imputation of urine-dipstick for albumin for the three validation  
588 cohorts does not adequately reflect UPCr. However, given that the level of proteinuria in chronic  
589 transplant nephropathy is generally mild, it is not expected that the imputation has major impact on  
590 the overall performance of the iBox score. During the discussion meetings (DMs) with the Applicant it  
591 was also evident that the qualification and validation exercise were tested separately for two different  
592 eGFR equations with equal performance (MDRD and SCr based CKD-EPI).

593 Model validation

594 Model diagnostics and Schoenfeld residual analysis for influential/outlier observations are adequate  
595 and do not cause concerns. The final model retained 8 variables in the full iBox score.

596 *Internal validation*

597 Internal validation focused on the full iBox panel. The abbreviated iBox Scoring System was not  
598 internally validated (p. 99 of the BD). The c-statistics for the derivation dataset were 0.809 and 0.803  
599 for the full and abbreviated iBox Scoring Systems, respectively (table 42, p. 100 BD). The c-statistics  
600 for the abbreviated iBox Scoring System showed that it is not significantly different than the c-  
601 statistics for the full iBox Scoring System. This is acknowledged and supports the use of both score  
602 variants.

603 Various scenarios and subpopulations were examined in the qualification dataset for their c-statistic  
604 using the iBox Scoring System (table 43, p. 102 BD). The full iBox Scoring System showed a good  
605 ability to discriminate the between higher and lower risk subjects for various important scenarios and  
606 subpopulations, with c-statistic values ranging from 0.76 to 0.87.

607 Three subsets showed significantly different c-statistic values from the c-statistic of 0.809 for the  
608 qualification derivation dataset (i.e., the 3,941 subjects for the full iBox Scoring System). This includes  
609 subjects transplanted with kidneys from elderly (c-statistic, 95% CI: 0.777, 0.746 to 0.808) and  
610 hypertensive donors (c-statistic, 95% CI: 0.771, 0.737 to 0.805). The proposed COU for the iBox  
611 Scoring System (i.e., evaluation at one-year post-transplant  $\pm$  28 days and censored at five-years and  
612 28 days post-transplant) in the mTORi subset of the derivation population shows also a good c-statistic  
613 value of 0.849 (95% CI from 0.804 to 0.893), suggesting the iBox Scoring System discriminates  
614 appropriately among subjects who meet the proposed COU. Overall, c-statistics in the derivation  
615 subsets suggests that the full iBox Scoring System performs well in various clinically relevant scenarios  
616 and subpopulations.

617 *External validation*

618 External validation was performed using the four external qualification datasets: Mayo Clinic Rochester  
619 and Helsinki University Hospital observational transplant center data, and the BENEFIT and BENEFIT-  
620 EXT RCTs. Analysis for these qualification validation datasets was restricted to the proposed COU, so  
621 only patients with full and abbreviated iBox Scoring System evaluations at one-year  $\pm$  28 days were  
622 retained for analysis, and data were censored at five-years and 28 days post-transplant.

623 The discrimination ability of the full and abbreviated iBox Scoring System models on each dataset was  
624 evaluated using the c-statistic censored at five-years plus 28 days post-transplant. All c-statistic  
625 values in table 45 are 0.70 or greater for each qualification validation dataset. The Applicant pointed  
626 out some shift in c-statistics score for the full- and the abbreviated iBox scores between datasets. This  
627 is explained by two participants with high eGFR at one year who lost their grafts. Similar change in c-  
628 statistic score was noted between the full- and abbreviated iBox in the Mayo cohort due to graft losses  
629 in two individuals at low risk of graft loss. The calibration results show that overall the predicted

630 number of events is reasonably well matching the number of observed events with some over- and  
631 underprediction when using single data sets, but with somewhat higher margins of error. This pertains  
632 to the full and abbreviated iBox score and also to subpopulations with treatment that is CNI based and  
633 without CNIs.

634 Overall, the data are considered encouraging. However, due to the limited number of data sets for  
635 validation and the limited number of graft loss events, the model assessment is subject to uncertainty  
636 and predicted event numbers show some variability. This also precluded assessment of the model in  
637 different subgroups, as was done for the derivation cohort.

#### 638 Supplementary analysis for validation

##### 639 *Competing risks analysis*

640 The sponsor used two methods for identifying whether the competing risk of death affects the full and  
641 abbreviated iBox Scoring System's predictions of graft loss. First, cumulative incidence functions (CIF)  
642 of graft loss that do and do not account for death were compared. Second, a Fine-Gray sub  
643 distribution survival model was built those accounts for death and compared to the iBox Scoring  
644 System, which is a Cox survival model that does not account for death. The sponsor gave detailed  
645 explanations and references for the two applied methods which are agreed upon.

646 The result of the analysis is that censoring deaths has little to no impact on predictions of graft loss in  
647 the derivation dataset.

##### 648 *Trial level surrogacy*

649 The focus of the TLS analysis was to: (1) estimate the treatment effect for each trial on full and  
650 abbreviated iBox Scoring System and graft loss, and (2) compute the correlation coefficient and/or the  
651 surrogate threshold effect (STE).

652 Due to limited availability of RCT, the two RCTs BENEFIT and BENEFIT-EXT were split into pseudo trials  
653 based on regions to support the TLS method. Splitting the data into pseudo-trials per region as  
654 performed by the Applicant (p. 121 BD) was helpful in allowing further assessment of the correlation  
655 at (pseudo-) trial level but does not contribute too much to improve precision of estimates  
656 for elucidating trial level surrogacy. Additionally, a subset of their derivation dataset was used,  
657 consisting of subjects who were on a CNI-free mTORi-based therapy, sirolimus or everolimus, versus  
658 CNI-based therapy at the time of transplant. To reduce potential confounding issues that can be  
659 present when examining non-RCT data propensity score techniques were used to reweight subjects in  
660 the derivation dataset. The addition of retrospective data from non-randomised comparisons in  
661 patients of the Loupy et al. cohort is only acceptable as supportive analysis. Inverse probability  
662 weighting based on propensity scores was used to allow comparisons. Results suggest that not all  
663 potential prognostic factors could be included, a stabilisation approach for the weights was necessary  
664 and it was not possible to generate bootstrap estimates for variance and correlation. While these  
665 issues raise concerns on the addition of non-randomised data to the exercise, even when these issues  
666 were not present, conclusions from the TLS analysis would not change.

667 No precise estimation of the trial-level correlation coefficient could be derived from these data. There  
668 are too few historical clinical trials available that are adequately sized and powered to quantitatively  
669 describe the treatment effect relationship on the surrogate and the true outcome. This prevented an  
670 adequate TLS analysis concerning whether the iBox Scoring System at one year detects a treatment  
671 effect that translates into differences in five-year death-censored allograft survival.

672 The low number of endpoint events in the available trials with an intervention limit establishing a  
673 correlation of treatment induced modification of the surrogate to treatment induced modification of the

674 endpoint at 5 years. Such a relation is considered key for establishing full surrogacy of a biomarker-  
675 based endpoint.

676 The TLS correlation analysis of treatment effects shows the limitations. The attempt to establish a trial  
677 level coefficient using a hierarchical Bayesian bivariate model shows a wide credible interval for the  
678 trial level coefficient including zero and therefore indicates the limitations for the precision of the  
679 estimate.

#### 680 Validation of an All-cause Endpoint score

##### 681 *ACE score development*

682 The primary event of interest in the ACE is all-cause allograft loss (including death). The abbreviated  
683 iBox composite score assessed at one year was used to assess all-cause graft survival. As can be  
684 expected, the model considerably underpredicts events. The model was therefore refined based on  
685 prior knowledge. Originally, known predictors for all-cause graft loss were delayed graft function (DGF)  
686 and rejection in the first year. These potential risk indicators were however not included in the model  
687 for predicting all-cause graft loss based on assessments at one-year post-transplant due to non-  
688 availability of rejection in the first year data in the PTG derivation dataset and due to “a non-  
689 substantial improvement” in risk prediction when DGF was included in the model compared to the use  
690 of the scoring system without DGF (= abbreviated iBox). The resulting model with eGFR, proteinuria  
691 and DSA was therefore taken forward for the ACE model. The considerations for model development  
692 are acknowledged.

693 With external validation datasets (p. 144 BD), C-statistics showed variable performance in moderate  
694 to good range of the ACE score on the discriminatory ability across the validation datasets (lowest C-  
695 statistic in Helsinki University Hospital 0.67 and Benefit-EXT RCT 0.67; C- statistic range from 0.67 to  
696 0.78). When excluding high risk donors, the model performed only slightly better (improvement in  
697 Helsinki University Hospital data plus 0.02, from 0.67 to 0.69).

698 C-statistics in the qualification derivation dataset (Loupy et al. 2019) showed moderate performance of  
699 the ACE score, again with a better performance when excluding high risk donors (C-statistics 0.75 with  
700 and 0.77 without high risk donors). Consequently, the model was adapted to exclude high risk donors.  
701 However, results showed moderate improvement in performance (table 82, p. 147 BD). The  
702 distribution plot of ACE scores for the derivation dataset without high-risk donors is illustrative (figure  
703 28, p. 147), showing separately the resulting counts for patients at 5 years discriminating patients  
704 alive with functional graft and deaths with functional graft and deaths with graft failure. This figure  
705 shows that the discriminatory ability for deaths with functional graft and deaths with graft failure of  
706 the ACE scoring system is modest.

707 The trial level surrogacy analysis (p. 149 BD) was repeated for the ACE. Treatment effect analyses  
708 were performed to investigate whether the treatment effect was significant on both the surrogate (ACE  
709 score) and the five-year all-cause graft survival. Two RCTs (Benefit-EXT RCT and BENEFIT RCT) and  
710 the mTORi derivation subsets were used. Concordance (significant treatment effect on ACE score and  
711 significant treatment effect on five-year all-cause survival) was found in one dataset (BENEFIT RCT),  
712 but not in the two others, where a directional effect on survival was found, but without statistical  
713 significance (likewise shown in analyses with and without high-risk donors). Like the surrogacy  
714 analysis in the iBox score systems, these results may be due to lack of statistical power.

715 Albeit not all predictors for all-cause graft survival were included in the ACE score, identity between  
716 this score and the abbreviated iBox score enables comparison of results. The performance of the ACE  
717 is less good than the iBox and this may be expected since predicting death events with functional graft  
718 may be difficult based on information tailored to predict renal events. Considering the above and the

719 observed results, from a performance and sensitivity perspective, the iBox score should be preferred  
720 over the ACE score.

#### 721 Conclusions

722 The Applicant initially proposed the iBox scoring system with full and abbreviated score without biopsy  
723 information as surrogate endpoint assessed at 1 year for prediction of outcomes at 5 years specifically  
724 tailored to settings with a conditional marketing authorisation application.

725 Overall, an extensive validation exercise has been performed, comprising internal validation based on  
726 prospectively collected data, external validation including randomised clinical studies and a trial-level  
727 surrogacy analysis. Previous work by Loupy et al. and the work done by the Applicant are  
728 comprehensively documented. Results show that the proposed iBox score models are suitable for  
729 individual predictions of graft loss events with good performance based on c-statistics and with the  
730 ability to predict numbers of graft loss events with reasonable, but not small margins of error.  
731 However, trial level surrogacy could not be established due to limited data in terms of available studies  
732 and event numbers. This is acknowledged by CHMP and also by the Applicant. Therefore, the iBox  
733 scores can currently not be formally qualified as surrogate endpoint for use as a primary endpoint.  
734 However, the use of the iBox as a secondary endpoint could be encouraged in order to further  
735 stimulate robust assessment of the iBox score and efficiency of drug development for treatments to  
736 prevent kidney graft failure.

737 During the discussion meetings with the Applicant it was evident that further data are needed in order  
738 to validate the iBox scoring system as a surrogate endpoint. This is understood and supported.

739

## 740 **4 Background as submitted by the applicant**

741 Please refer to the validated Briefing Document (BD) published as separate document for the evidence  
742 presented.