

1 **Predicting Success of Phase III Trials in Oncology**

2 Meta-analysis of phase II and phase III trials initiated between 2003 and 2012 to generate a  
3 decision analytical model predicting success of future phase III trials

4

5 Stephan Hegge<sup>1,#</sup>, PhD, Markus Thunecke<sup>2</sup>, PhD, Matthias Krings<sup>2</sup>, PhD, Léonard Ruedin<sup>3</sup>,  
6 MSc, Jan Saputra Müller<sup>4</sup>, MSc, Paul von Büнау<sup>3</sup>, PhD

7

8 <sup>1</sup>Hegge Holding UG, Kopernikusstr. 24, 10247 Berlin

9 <sup>2</sup>Catenion GmbH, Münzstrasse 18, 10278 Berlin, Germany

10 <sup>3</sup>idalab GmbH, Potsdamer Straße 68, 10785 Berlin, Germany

11 <sup>4</sup>AskBy GmbH, Boxhagener Str. 18, 10245 Berlin

12

13 #corresponding author

14 Email: [StephanHegge@gmx.de](mailto:StephanHegge@gmx.de)

15 **Word Count: 3493**

## 16 **Key Points**

17 **Question:** What is the probability of success (PoS) for single phase III (PhIII) trials in  
18 oncology?

19 **Findings:** We developed a model allowing to predict the PoS of single PhIII trials in oncology  
20 with a predictive performance of 73%<sub>PP</sub> and demonstrated that qualitative factors such as  
21 strength of PhII knowledge and sponsor R&D strength can be captured in quantitative scores  
22 that have significant predictive power.

23 **Meaning:** The model can help study sponsors to analyze and amend planned clinical trials,  
24 and investors to choose where to invest best.

25

## 26 **Abstract**

27 **Importance:** We developed a model predicting the probability of success (PoS) for single  
28 planned or ongoing PhIII trials based on information available at trial initiation. Such a model  
29 is highly relevant for study sponsors to capture risk and opportunity on a trial-to-trial basis  
30 through trial optimization, and for investors to select drugs whose trial design match their  
31 investment strategy.

32 **Objectives:** To predict the outcome of planned or ongoing PhIII trials in oncology, given  
33 publicly available prior information

34 **Design, Setting, Participants:** Predictive modeling using publicly available data for 360  
35 completed PhIII and 1240 PhII studies initiated between 2003 and 2012. Success and failure  
36 of PhIII studies were modeled using Bayesian logistic regression model.

37 **Main Outcome Measures:** Predicted PoS of individual PhIII trials based on a Bayesian  
38 model calibrated on publicly available data translated into 16 composite scores. Those  
39 scores cover aspects such as trial design, indication, number of patients, phase II (PhII)  
40 study outcomes, experience of sponsor at time of trial initiation, and others.

41 **Results:** The model allows to calculate the PoS distribution – including credible intervals –  
42 for a PhIII trial in oncology. The predictive performance was determined using an area under  
43 the receiver-operator curve (AUROC), resulting in an overall performance of 73%<sub>oPP</sub> (mean  
44 AUROC). We identified two key factors contributing to the predictive performance of the  
45 model: quality and strength of PhII data and experience of the sponsor at the time of study  
46 initiation.

47 **Conclusion and Relevance:** We describe the generation and application of a statistical  
48 model predicting the PoS for individual PhIII trials in oncological indications with  
49 unprecedented predictive performance. Compared to other approaches, this is the first study  
50 generating a fully transparent model resulting in trial-specific PoS distributions. Moreover, we  
51 have shown that qualitative concepts such as PhII knowledge or sponsor R&D strength can  
52 be captured in quantitative scores and that these scores have a high predictive power.

53

## 54 **Introduction**

55 In recent years, predictive algorithms have become ubiquitous across a wide range of  
56 industries, such as logistics – e.g. Amazon’s predictive shipping<sup>1</sup> – or information retrieval –  
57 e.g. Google’s predictive search algorithm<sup>2</sup>. By combining information from a vast number of  
58 sources in an objective, unbiased manner, predictive algorithms can outperform human  
59 decision making with respect to accuracy and speed at marginal cost. Even in the public  
60 sector, political decision makers have become increasingly aware of the importance of  
61 accurate predictions and have started evaluating different approaches in forecasting  
62 tournaments<sup>3</sup>. Predictive algorithms can aid decision makers in the pharmaceutical industry  
63 as well. However, so far, adoption has been limited. The current decision making process,  
64 from discovery to clinical development phases, is characterized by a series of decision points  
65 defined by formal go/no-go criteria<sup>4</sup> that relate to the available clinical data. This process is  
66 implemented to aid executives as it reflects regulatory requirements for each indication.  
67 However, it has been demonstrated that decisions based on real world cases vary greatly,  
68 ranging from absolutely ‘go’ to absolutely ‘no-go’ due to subjective interpretation of identical  
69 data<sup>5</sup>.

70 We argue that predictive algorithms can be employed for improving and rationalizing decision  
71 making in Pharma, particularly in clinical trials representing the most crucial and expensive  
72 centerpiece of drug development. Predictive algorithms based on Big Data have  
73 demonstrated that they are able to aid or even outperform drug developers and physicians  
74 when it comes to predicting either patient accrual rates in clinical trials<sup>6,7</sup>, or optimal cancer  
75 rehabilitation<sup>8</sup>, or supportive care interventions<sup>8</sup>.

76 Trial decision making for any drug in clinical development has far-reaching implications. On  
77 one hand, hundreds to thousands of patients are recruited to test the effects, each of them  
78 hoping to benefit from the new drug. On the other hand, sponsors risk tens to hundreds of  
79 million Dollars on trials that may or may not demonstrate the drug’s effect<sup>9</sup>. It is therefore in  
80 the interest of sponsors to terminate failures early, without compromising on the quality of

81 clinical trials or terminating successful agents, and – most importantly – without exposing  
82 clinical trial subjects to unnecessary risk<sup>10</sup>. It is also in the interest of patients to participate in  
83 trials that have a strong chance to achieve a positive outcome.

84 Success and failure rates in pharmaceutical drug development are described by several  
85 sources, ranging from commercial benchmarking providers like BioMedTracker<sup>11,12</sup> to  
86 academic papers<sup>13–16</sup>.

87 Still, no established method exists to predict the probability of success (PoS) for an individual  
88 clinical trial; in fact, there is no clear definition for trial success and failure. Hence current  
89 approaches<sup>11–15</sup> focus on determining the average attrition rates from phase to phase – so  
90 called phase transitions counts (PTC) – during drug development resulting in historical  
91 success rates (hSRs). Phase transitions – and hence hSRs – are either assessed on an  
92 indication<sup>11</sup> (hSR<sub>Indication</sub>; eFig 1, eTable 1), or on a program level<sup>10,17</sup> (hSR<sub>Program</sub>, eFig 1,  
93 eTable 1). It has become gold standard in trial planning and portfolio management to take  
94 hSRs and use them as forward-looking estimates of PoS. Moreover, using hSR as gauge for  
95 forward-looking PoS is limited to factors that allow for straightforward stratification. Drivers  
96 that are more complex in nature, such as the knowledge a sponsor has accumulated from  
97 prior phases, cannot be analysed using simple hSRs. This approach has far-reaching  
98 implications in decision making for several reasons:

- 99 a) hSRs are historic, past transition rates, however they are naively used as forward-  
100 looking probabilities, without a statistical evaluation of their predictive power
- 101 b) hSRs are averages. Thus, applying these as forward-looking PoS to specific  
102 individual trials will systematically overestimate the PoS of riskier than average  
103 drug development programs, while underappreciating the PoS of particularly well  
104 conducted low-risk programs
- 105 c) hSR-based PoS values are typically used as point estimates without credible  
106 intervals<sup>11</sup>. This can lead to a false sense of confidence, and falsely informed  
107 decisions – such as terminating a drug that would work and be safe in favor of

108 continuing a drug that is ineffective or unsafe – when rating the PoS of one  
109 program higher than another

110 d) Neither of the hSR approaches consider individual trials, but are limited to  
111 indications<sup>11</sup>, drug programs<sup>10</sup> or even entire therapeutic areas<sup>17</sup>

112 Here we present an approach providing forward looking PoS estimates for individual PhIII  
113 trials in oncology based on publicly available data. Firstly, all potentially relevant parameters  
114 were classified and translated into numerical or categorical data. Secondly, we identified and  
115 quantified correlations between each parameter and trial outcome. Thirdly, we developed  
116 and calibrated a Bayesian algorithm for predicting trial outcomes, which provides a full  
117 posterior distribution. Fourthly, we have shown that complex qualitative factors such as PhII  
118 knowledge or sponsor R&D strength can be modeled using quantitative scores and that  
119 these scores are indeed significantly predictive. This approach allows decision-makers to  
120 capture risk and opportunity on a trial-to-trial basis, and helps trial sponsors and planners to  
121 identify and mitigate trial- and indication-specific risks.

122

## 123 **Materials and Methods**

### 124 **Data Sources**

125 ClinicalTrials.gov (CT.gov) is a publicly accessible database run by the United States  
126 National Library of Medicine at the National Institutes of Health<sup>18</sup>. It served as a starting point  
127 for our analysis as it is the largest registry for clinical trials to date providing general  
128 characteristics and outcome information of about ~283,000 (as of January 2018) federally  
129 and privately supported clinical trials. 47% of studies recruit patients outside the US, 36% of  
130 trials are recruiting in the US only, and 6% recruit patients in both, US and non-US locations,  
131 while 11% do not provide that information<sup>18</sup>.

132 To complement for potentially publicly available information not listed on CT.gov we used the  
133 European Clinical Trial Database<sup>19</sup>, JAPIC Clinical Trials Information<sup>20</sup>, PubMed<sup>21</sup> and Adis  
134 R&D<sup>22</sup> (ADIS). We used PubMed to identify publications associated with each given trial.  
135 ADIS is a database for drug development gathering all available information (*in vitro* and *in*  
136 *vivo* experiments, preclinical data, clinical trial outcomes, chemistry, conference posters,  
137 conference abstracts, press releases etc.) for a given compound.

### 138 **Study Sample**

139 We queried the database of CT.gov<sup>18</sup> (>183,000 trials by April 2014) and applied filters to  
140 identify a sample of novel therapeutics (i.e., new chemical entities (NCEs) or novel biological  
141 entities (NBEs)) for which a PhIII study in oncology had been initiated between 01.01.2003  
142 and 01.03.2012 (Fig 1). Filter criteria excluded generics, biosimilars, reformulations,  
143 radiotherapy, and combination therapies of two or more non-novel therapeutics. Cell-  
144 therapies, non-therapeutic agents, such as diagnostic tools, contrast agents and supportive  
145 care approaches were also not considered. We removed all duplicated records.

### 146 **Technical Implementation**

147 For MCMC sampling, we used the software package JAGS version 4.3. All other  
148 computations were performed in Python version 2.7, relying heavily on pandas version 0.17.

## 149 **Results**

### 150 *Database*

151 ClinicalTrials.gov is the largest registry of clinical trials database with >183,000 registered  
152 trials at the time of query (Fig 1, I). Strict application of exclusion criteria resulted in 360 PhIII  
153 trials across 37 oncology indications (Fig 1, IV). The 111 NBEs and NCEs associated with  
154 these PhIII studies served as a starting point for the complementary analysis, which aimed to  
155 gather all publicly available PhII and PhIII information associated with the original set of 360  
156 PhIII studies by means of searching for the drug name, the trial identifier, trial names and its  
157 synonyms. We searched in the European and Japanese trial registries EudraCT and JAPIC,  
158 scanned published literature using PubMed.gov, company homepages, and checked Adis  
159 R&D, a database for drug research and development (Fig 1, V). This approach gave rise to  
160 1240 PhII studies, 488 publications, and 111 ADIS data entries, respectively, with detailed  
161 information about PhII and PhIII study outcomes, endpoint measures, study design, actual  
162 patient population etc. (eTables 1-3). Please note, PhI results were not considered for this  
163 proof-of-principle analysis, given PhI studies are less well defined with regard to patient  
164 stratification (mixed patient population with regard to tumor type, tumor stage, line of  
165 treatment) and endpoint measures. We defined trial success by meeting all primary  
166 endpoints (see eTable 3 for detailed classification).

### 167 *Scoring approach for predictive modeling*

168 The key step in predictive modeling is to construct informative (predictive) scores from the  
169 available raw data. To this end, we employed a hybrid approach that combines (a) expert-  
170 driven design of complex scores to incorporate human judgement and experience and (b) a  
171 purely data-driven approach to assigning weights to variables and then selecting the most  
172 predictive variables. The group of domain experts consisted of eight consultants with an  
173 academic life science background and several years of work experience in pharma R&D.  
174 Notably, expert input was not employed with regard to individual trials (which would lead to  
175 biased results), but exclusively in the structural design of the scoring methodology. Fig 2

176 provides a conceptual overview of this approach illustrated by an exemplary drug X, which is  
177 developed in three different indications A, B, and C (Fig 2A). We identified a battery of  
178 information available at a given point in time (Fig 2A, time of assessment, vertical grey line)  
179 to assess the probability of success for a given PhIII trial of interest (TOI, blue box). We  
180 classified all available information into time-dependent variables (Fig 2B, eTable 1), drug-  
181 related characteristics (Fig 2C, eTable 2) and trial-specific characteristics (Fig 2D, eTable 3).  
182 Within each of these categories we created composite scores combining weight and  
183 magnitude of information. All composites were broken down to objective and quantifiable  
184 elements (Fig 2B i-iii). We developed unbiased decision matrices for variables that offered no  
185 direct read-out from the database (e.g. Relatedness between PhII and PhIII TOI, eFig 2). In  
186 total, we created three classes of parameters comprising 16 composite scores (see eTables  
187 2-4) based on 82 variables. Consequently, each PhIII trial is described by a unique set of  
188 descriptors, allowing the generation of a trial-specific prediction.

### 189 *Generation of predictive model*

190 The overall approach to predictive modeling is a dynamic Bayesian logistic regression<sup>23</sup>.  
191 More specifically, the log-odds of the binary target variable (Fig 3, Outcome of TOI S either  
192 success or failure) is modeled as a linear function of the independent variables with a  
193 Gaussian prior for each variable  $k$ 's weight  $\alpha$  (Fig 3A) resulting in a trial specific PoS  $\Theta$ . The  
194 coefficients (weights) are estimated in a Bayesian approach using a Markov Chain Monte  
195 Carlo (MCMC)<sup>24</sup> simulation to generate samples from the posterior of each parameter (Fig  
196 3B); the final estimate is the posterior mean (Fig 3C).

### 197 *Training and evaluation of the predictive model*

198 As a performance metric, we use the Area under the Receiver-Operator Characteristic  
199 (AUROC)<sup>25</sup>. To account for the time structure in the data (PhIII trial initiation spanning 2003-  
200 2012), we use a dynamic modeling approach with regard to (a) the construction of the  
201 independent variables from raw data (Fig 4A), (b) variable selection (Fig 4B) and (c) overall  
202 performance evaluation (Fig 4C). Regarding step (a), this means that in the computation of

203 scores only such information is used that was available at the time point of the respective  
204 PhIII trial's initiation date (eTable 1). For variable selection (b) and overall performance  
205 evaluation (c) we adopt a time-series cross-validation strategy.

#### 206 *(Ir-)Relevance of variables*

207 The predictive performance of single composite scores can be calculated for 15 out of the 16  
208 composite scores that we designed, i.e. all scores but novelty of MoA (Fig 4A). When  
209 applying the AUROC performance metric, we found positive correlation – mean AUROC  
210 >55%, ranging up to 64% AUROC – with success or failure for the metrics SPIIK, Company  
211 R&D strength, number of subjects in trial, designations, trial type and prior registrations of the  
212 drug in other indications.

213 We found no correlation – AUROC >45% and <55% - with success or failure for the metrics  
214 modality, indication, geography, and involvement of Big Pharma.

#### 215 *Model selection*

216 The best full predictive risk model (PRM) was built based on the single composite scores  
217 with a positive contribution to the model, while scores without impact were omitted (Fig 4B).  
218 To arrive at the best PRM, we added – at each step – the variable producing the highest  
219 mean AUROC increase of the model (or the smallest decrease). The variables selected (Fig  
220 4B) are those that bring the highest overall predictive performance, reported as the mean  
221 AUROC over five time-series cross-validation splits (Fig 4C, Fig 3S). Note that the predictive  
222 performance of the single scores is not strictly additive due to overlapping information.

#### 223 *Overall model performance*

224 The performance of the model with the best combination is ~73%<sub>OPP</sub> and includes 12  
225 variables (Fig 4C, eFig 3B, dotted line). In other words, confronted with one successful and  
226 one unsuccessful trial, using the PRM one will correctly identify the successful trial in ~73%  
227 of the cases by picking the trial with the higher predicted PoS.

#### 228 *Exemplary model outcomes*

229 To illustrate the model's output on single trials, we elected to highlight the variables  
230 characterizing exemplary clinical trials and to visually signify their influence on the PoS  
231 prediction (Fig 5A). Consequently, we selected two clinical studies at different ends of the  
232 spectrum, one with low mean PoS – gefitinib in treating patients with esophageal cancer that  
233 is progressing after chemotherapy, NCT01243398<sup>26</sup>, Fig 5B – and one with high mean PoS –  
234 trastuzumab emtansine versus capecitabine + lapatinib in participants with HER2-positive  
235 locally advanced or metastatic breast cancer (EMILIA), NCT00829166<sup>27</sup>, Fig 5C.

236 The initiation dates of both studies were respectively March 2009<sup>26</sup> for NCT01243398 and  
237 February 2009<sup>27</sup> for NCT00829166. To mimic a real world PhIII decision point, the model  
238 goes back to the time before those particular PhIII were initiated, thereby only accessing data  
239 that were available then. The resulting PoS distributions are therefore built on 92 PhIII  
240 studies and its corresponding PhII trials preceding initiation of NCT01243398 and  
241 NCT00829166.

242 According to our model, the PhIII trial in in esophageal cancer with gefitinib had a predicted  
243 mean  $PoS_{\text{Trial}}$  of 7.8%<sub>PoS</sub> with 95%<sub>CI</sub> credible interval ranging from  $PoS_{\text{Trial}} = 0.3\%$ <sub>PoS</sub> to  
244 67.2%<sub>PoS</sub> at time of initiation. The study eventually failed to meet the primary endpoint  
245 median overall survival<sup>28</sup> (Fig 5D). On the other hand, according to our model the PhIII trial in  
246 breast cancer with trastuzumab emtansine had a predicted mean  $PoS_{\text{Trial}}$  of 88.6%<sub>PoS</sub> with  
247 95%<sub>CI</sub> credible interval ranging from  $PoS_{\text{Trial}} = 19.6\%$ <sub>PoS</sub> to 99.6%<sub>PoS</sub> at time of initiation. That  
248 clinical study was successful<sup>29</sup> (Fig 5E), confirming the predicted high probability of success.

## 249 Discussion

250 Drug developers use historical success rates as forward-looking estimates to determine the  
251 PoS of an individual trial to inform their decision making. These PoS rates are often adjusted  
252 based on the opinion of subject matter experts, so-called Key Opinion Leaders (KOLs). For  
253 more than 60 years studies<sup>30-32</sup> have kept demonstrating that expert opinions are no better  
254 than guessing, while in the last decade or so it was established that algorithms are able to  
255 aid or even outperform drug developers and physicians when it comes to predicting patient  
256 accrual rates in clinical trials<sup>33,34</sup>, optimal cancer rehabilitation, or supportive care  
257 interventions<sup>8</sup>. Therefore, it is surprising that an industry praising rationale drug development  
258 still takes decisions regarding investments, strategy, and science based on subjective expert  
259 advice. While the traditional approach certainly has some merit, we argue that a data-driven  
260 prediction can complement – if not quite replace – traditional methods such as KOL  
261 interviews and hSR benchmarking. In particular, we have demonstrated that the  
262 consideration of complex drivers of PhIII success such as the accumulated knowledge from  
263 prior phases is not limited to the judgement of experts (KOLs), but can also be addressed in  
264 an empirical data-driven manner using a sophisticated scoring approach presented in this  
265 study.

### 266 *Discussion of results*

267 Employing several publicly available databases<sup>18,21,22</sup> we developed a predictive model  
268 generating forward-looking and trial-specific probability of success (PoS) distributions for  
269 PhIII trials in oncology.

270 The most relevant single composite scores contributing to the full PRM are Strength of Ph II  
271 Knowledge (SPIIK) (mean AUROC<sub>SPIIK</sub> = 64%) and Company R&D Strength at time of trial  
272 initiation (mean AUROC<sub>SPIIK</sub> = 63%, Fig. 4A).

273 The SPIIK includes the strength and relatedness of the combined PhII evidence that exists  
274 before starting the PhIII (Fig 2, eTable 2). Our use of a decision tree optimized for prediction  
275 allowed us to model the relevance of a combined PhII body of evidence to a particular PhIII

276 study and its design (eFig 2). The predictive performance of SPIIK alone confirms that the  
277 information building this composite score is relevant indeed. This is in line with both, common  
278 sense and regulatory guidance<sup>35</sup>, suggesting that the sum of PhII results are to some degree  
279 indicative for PhIII trial outcomes.

280 The sponsor's past track record in oncology (Company R&D Strength) had the second  
281 largest predictive performance for PhIII studies. As this criterion includes both the number of  
282 past PhII and III studies in oncology as well as the outcomes, it essentially describes where  
283 an organization stands on the learning curve when it comes to designing studies in oncology.  
284 Noteworthy, this value is time-dependent to consider the situation on the day of trial initiation.  
285 There is a strong effect of having designed phase II and III studies that met their primary  
286 endpoints on the ability to do it again. Janssen, Celgene, and Genentech are the top 3  
287 performers in this category of companies with at least 10 PhIII studies (2003 – 2012) in our  
288 sample.

#### 289 *The role of indication*

290 The factor 'indication' provides no additional value to the predictive risk model (Fig 4A, B).  
291 Notably, this is in line with expectations due to the technical bias in trial selection (Fig 1); We  
292 started our search with PhIII trials and only subsequently enriched with PhII trials associated  
293 with the selected PhIII trials. Therefore, we introduced a bias for drugs that made it into PhIII  
294 in at least one indication. Drugs exclusively developed in indications known for high failure  
295 rates (e.g. Pancreatic Cancer)<sup>14</sup> hardly make it into a PhIII trial in the respective indication,  
296 hence PhIII trials in these indications are underrepresented in this proof of concept study.

#### 297 *Novelty of MoA*

298 In order to factor in the degree of innovation brought about by the compounds investigated in  
299 PhIII, we designed a composite score taking into account the novelty of MoA. That composite  
300 score was excluded from the model, as the results of our attempt were not conclusive. On  
301 one hand, this is due to the complexity of embodying the qualitative nature of innovation into  
302 a quantitative variable. On the other hand, there is a lack of availability of systematic,

303 comprehensive data due to fundamental differences in the MoA classification schemes used  
304 and the level of information provided by companies (eTables 2 and 4).

### 305 *Comparison to other approaches*

306 Empirical studies of clinical trial success broadly fall into two categories: (i) retrospective  
307 descriptive analysis of success rates<sup>11,12,15,16,36,37</sup> and (ii) predictive approaches to modeling  
308 clinical trial success<sup>10,38</sup>, as in the present study. From a statistical methodology point of  
309 view, retrospective descriptive studies focus on estimates of success rates computed from  
310 empirical binary (fail vs. success)h<sup>11,39</sup>. Confidence intervals for the reported PoS estimates  
311 are mostly not provided, with the exception of Wong 2018<sup>16</sup> (standard errors).

312 Among the predictive studies, a much wider range of methodological approaches can be  
313 found in the literature (eTable 5). Schachter et al.<sup>10</sup> employed a Bayesian Network, a highly  
314 flexible model class, which could potentially be used to formulate expressive bottom-up  
315 generative models of trial success. Still, the approach was limited at the time due to scarcity  
316 of available databases and lack of historical data<sup>40</sup>, resulting in a hold-out validation set too  
317 small (n=14) to generate reliable outcomes. DiMasi *et al.*<sup>38</sup> employ an unorthodox type of  
318 predictive model which uses a scoring logic to compute the predicted PoS for a given trial.

319 In contrast to others, we use a linear regression model to reduce overfitting and for ease of  
320 interpretation but calibrate parameters in a Bayesian fashion so that credible intervals for  
321 parameter estimates and a posterior predictive distribution is available for PoS estimates.  
322 This is the basis for downstream Monte-Carlo simulation of portfolio-level effects. Secondly,  
323 for model evaluation we use a time-series cross-validation strategy to analyze the  
324 performance over time, as more historic information becomes available. This analysis shows  
325 not only the mean AUROC on one data set but provides also variation and stability of  
326 predictive performance.

### 327 *Limitations of approach*

328 The current algorithm is focused on oncology PhIII trials. For this proof of concept study, we  
329 chose oncology over other therapeutic areas, because trial endpoints (mPFS and mOS)  
330 across all oncology indications are both, quantifiable and comparable in nature, providing a  
331 strong foundation for modeling approaches.

332 We excluded several non-standard modalities including the cellular therapies which are  
333 changing the treatment landscape as we speak. In principle, the algorithm is also prone to  
334 certain regulatory aspects (e.g. break through designations) that allow a development  
335 program to move from phase I (PhI) directly to PhIII or approval, respectively.

### 336 ***Next steps & outlook***

337 Based upon this proof of concept study, the model can potentially be expanded to (1) predict  
338 PhII trial outcomes based on data from pre-clinical and PhI studies, (2) therapeutic areas  
339 other than oncology (e.g. cardiovascular diseases), (3) incorporate more modalities (e.g.  
340 CAR T cells) for which a growing body of evidence is becoming available, and to (4) allow for  
341 the integration of non-public information available to drug developers (sponsors and  
342 investigators) in cooperation with the project teams.

### 343 ***Conclusions***

344 The algorithm presented here can distinguish successful from unsuccessful trials with much  
345 greater confidence than any other publicly available approach reviewed<sup>10,12,14–17,38</sup>. The  
346 positive predictive value can be tuned up to >80%<sub>oPPV</sub> by accepting more false negatives  
347 (lower sensitivity). To our knowledge, this is the first approach allowing to quantitatively  
348 predict the probability of success for single trials. Our model uses publicly available  
349 information only, including that of prior trials with perhaps only remote relatedness to the trial  
350 in question, and then delivers a specific prediction for a given trial. In addition, the model is  
351 fully transparent, adaptive on a trial-to-trial basis, provides unprecedented granularity (e.g.  
352 consideration of line of treatment, or background therapy) and allows identification of factors  
353 negatively (or positively) influencing the trial's predicted PoS<sub>Trial</sub>.

354 Such an algorithm has a number of obvious applications of high medical, strategic and  
355 financial value, quite apart from the ethical dimension of a doctor's decision to enroll patients  
356 in a study. Both sponsors and investors involved in the field of oncology could benefit greatly  
357 from a predictive algorithm assessing the prospects of a specific study, in particular by

- 358 • Supporting sponsoring companies to maximize success by designing their individual  
359 studies based on the highest possible  $Pos_{Trial}$
- 360 • Helping investors determine the impact of PhIII outcomes on valuation. This is especially  
361 relevant for those biotechs with a single PhIII asset. In addition, investors able to pursue  
362 different strategies could identify trials (and companies behind the studies) that match  
363 their investment strategy, e.g. pick-the-winner-drop-the-loser or vice versa.

## 364 **Acknowledgments**

365 SH drafted manuscript and figures, developed concept, acquired, analyzed and interpreted  
366 data, managed project. MT and MK supervised project, challenged concept, provided  
367 material and technical support, and critically reviewed manuscript. LR drafted manuscript and  
368 figures, acquired, analyzed and interpreted data. JSM developed, programmed and  
369 calibrated the model, retrieved data and critically reviewed manuscript. PvB developed the  
370 model, provided statistical analysis, interpreted data, drafted the manuscript, and supervised  
371 the project.

372 The authors thank Birte Arlt, Nina Heid and Jennifer Price for data curation and classification,  
373 Moritz Neeb and Daniel Kirsch for advice on data architecture and for co-developing of  
374 algorithms required for model, Gerrit Buurmann for strategic advice and for expert  
375 classification of trial data, and Johannes Zimmermann for critical review of manuscript.

376 Conflict of Interest: No conflicts of interest are declared for JSM, PvB, MK. SH, LR and MT  
377 have personal investments in several biotechnology companies. SH is Senior Director of  
378 Corporate Strategy at HotSpot Therapeutics Inc. and general manager of Hegge Holding  
379 UG. JSM is co-founder of AskBy GmbH. No funding bodies had any role in study design,  
380 data collection and analysis, decision to publish, or preparation of the manuscript. The  
381 authors were personally salaried by their institutions during the period of writing (though no  
382 specific salary was set aside or given for the writing of this paper).

## 383 References

- 384 1. Spiegel JR, McKenna MT, Lakshman GS, Nordstrom PG. *Method and System for*  
385 *Anticipatory Package Shipping*. Google Patents; 2013.  
386 <https://www.google.com/patents/US8615473>.
- 387 2. Hu W, Zhang X, Bolivar A, Shoup RS. *Predictive Algorithm for Search Box Auto-*  
388 *Complete*. Google Patents; 2015. <http://www.google.com.pg/patents/US20150193449>.
- 389 3. Mellers B, Stone E, Murray T, et al. Identifying and cultivating superforecasters as a  
390 method of improving probabilistic predictions. *Perspect Psychol Sci*. 2015;10(3):267–281.
- 391 4. Hedner T, Cowlrick I, Wolf R, Olausson M, Klofsten M. The changing structure of the  
392 pharmaceutical industry: perceptions on entrepreneurship and openness. 2011.
- 393 5. Cowlrick I, Hedner T, Wolf R, Olausson M, Klofsten M. Decision-making in the  
394 pharmaceutical industry: analysis of entrepreneurial risk and attitude using uncertain  
395 information. *RD Manag*. 2011;41(4):321-336.
- 396 6. London JW, Balestrucci L, Chatterjee D, Zhan T. Design-phase prediction of potential  
397 cancer clinical trial accrual success using a research data mart. *J Am Med Inform Assoc*.  
398 2013;20(e2):e260-e266. doi:10.1136/amiajnl-2013-001846
- 399 7. Cheng SK, Dietrich MS, Dilts DM. Predicting Accrual Achievement: Monitoring  
400 Accrual Milestones of NCI-CTEP-Sponsored Clinical Trials. *Clin Cancer Res*.  
401 2011;17(7):1947-1955. doi:10.1158/1078-0432.CCR-10-1730
- 402 8. Buffart LM, Kalter J, Chinapaw MJ, et al. Predicting Optimal Cancer Rehabilitation  
403 and Supportive care (POLARIS): rationale and design for meta-analyses of individual patient  
404 data of randomized controlled trials that evaluate the effect of physical activity and  
405 psychosocial interventions on health-related quality of life in cancer survivors. *Syst Rev*.  
406 2013;2(1):75.
- 407 9. Speich B, von Niederhäusern B, Schur N, et al. Systematic review on costs and  
408 resource use of randomised clinical trials shows a lack of transparent and comprehensive  
409 data. *J Clin Epidemiol*. December 2017. doi:10.1016/j.jclinepi.2017.12.018
- 410 10. Schachter AD, Ramoni MF, Baio G, Roberts TG, Finkelstein SN. Economic  
411 Evaluation of a Bayesian Model to Predict Late-Phase Success of New Chemical Entities.  
412 *Value Health*. 2007;10(5):377-385. doi:10.1111/j.1524-4733.2007.00191.x
- 413 11. Hay M, Thomas DW, Craighead JL, Economides C, Rosenthal J. Clinical  
414 development success rates for investigational drugs. *Nat Biotechnol*. 2014;32(1):40–51.
- 415 12. Thomas DW, Burns J, Audette J, Carroll A, Dow-Hygelund C, Hay M. *Clinical*  
416 *Development Success Rates 2006-2015*. BIO Industry Analysis. Amplion, Inc.,  
417 Biomedtracker, Biotechnology Innovation Organization (BIO); 2016:1-26.
- 418 13. DiMasi JA, Grabowski HG. Economics of New Oncology Drug Development. *J Clin*  
419 *Oncol*. 2007;25(2):209-216. doi:10.1200/JCO.2006.09.0803
- 420 14. DiMasi JA, Reichert JM, Feldman L, Malins A. Clinical Approval Success Rates for  
421 Investigational Cancer Drugs. *Clin Pharmacol Ther*. 2013;94(3):329-335.  
422 doi:10.1038/clpt.2013.117
- 423 15. DiMasi JA. Pharmaceutical R&D performance by firm size: approval success rates  
424 and economic returns. *Am J Ther*. 2014;21(1):26–34.

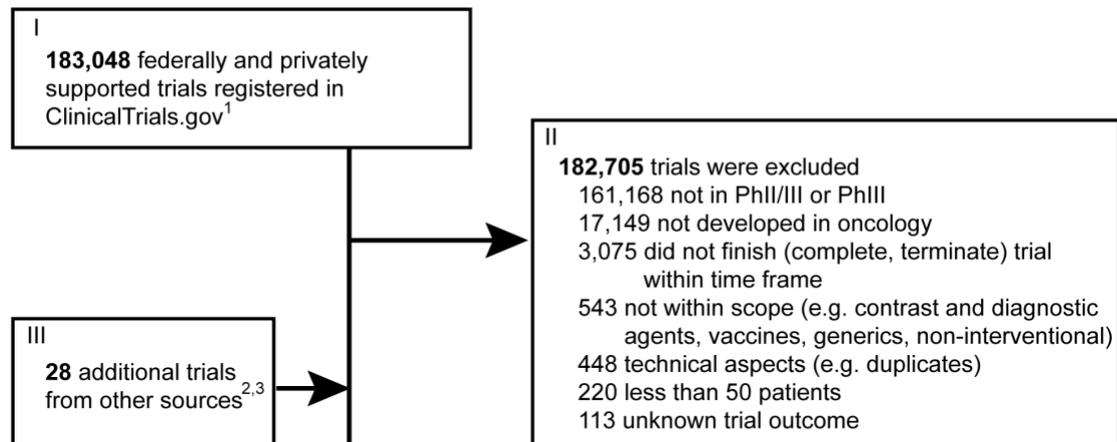
- 425 16. Wong CH, Siah KW, Lo AW. Estimation of clinical trial success rates and related  
426 parameters. *Biostatistics*. January 2018. doi:10.1093/biostatistics/kxx069
- 427 17. *CMR International 2015 Pharmaceutical R&D Factbook*. Thomsen Reuters, Centre  
428 for Medicines Research (CMR) International; 2015.
- 429 18. Trends, Charts, and Maps - ClinicalTrials.gov.  
430 <https://clinicaltrials.gov/ct2/resources/trends>. Accessed April 4, 2014.
- 431 19. EudraCT. European Clinical Trial Database. [https://eudract.ema.europa.eu/results-](https://eudract.ema.europa.eu/results-web/index.xhtml)  
432 [web/index.xhtml](https://eudract.ema.europa.eu/results-web/index.xhtml). Accessed December 4, 2014.
- 433 20. JAPIC Clinical Trials Information. JAPIC Clinical Trials Information.  
434 [http://www.clinicaltrials.jp/user/cteSearch.jsp;jsessionid=33025EAEEDA366B0FA4A26F45B](http://www.clinicaltrials.jp/user/cteSearch.jsp;jsessionid=33025EAEEDA366B0FA4A26F45B1F15DF)  
435 [1F15DF](http://www.clinicaltrials.jp/user/cteSearch.jsp;jsessionid=33025EAEEDA366B0FA4A26F45B1F15DF). Accessed December 4, 2014.
- 436 21. pubmeddev. Home - PubMed - NCBI. <https://www.ncbi.nlm.nih.gov/pubmed/>.  
437 Accessed December 4, 2014.
- 438 22. Adis R&D. [www.springer.com](http://www.springer.com). <https://www.springer.com/gp/librarians/adis-r-d/7790>.  
439 Accessed December 4, 2014.
- 440 23. Gelman A, Carlin JB, Stern HS, Dunson DB, Vehtari A, Rubin DB. *Bayesian Data*  
441 *Analysis, Third Edition*. CRC Press; 2013.
- 442 24. W.K. Hastings. Monte Carlo Sampling Methods Using Markov Chains and Their  
443 Applications. *Biometrika*. 1970;57(1):97-109.
- 444 25. Zweig MH, Campbell G. Receiver-operating characteristic (ROC) plots: a fundamental  
445 evaluation tool in clinical medicine. *Clin Chem*. 1993;39(4):561-577.
- 446 26. Gefitinib in Treating Patients With Esophageal Cancer That is Progressing After  
447 Chemotherapy - Tabular View - ClinicalTrials.gov.  
448 <https://clinicaltrials.gov/ct2/show/record/NCT01243398>. Accessed November 17, 2018.
- 449 27. A Study of Trastuzumab Emtansine Versus Capecitabine + Lapatinib in Participants  
450 With HER2-positive Locally Advanced or Metastatic Breast Cancer - Tabular View -  
451 ClinicalTrials.gov. <https://clinicaltrials.gov/ct2/show/record/NCT00829166>. Accessed  
452 November 17, 2018.
- 453 28. Petty RD, Dahle-Smith A, Stevenson DAJ, et al. Gefitinib and *EGFR* Gene Copy  
454 Number Aberrations in Esophageal Cancer. *J Clin Oncol*. 2017;35(20):2279-2287.  
455 doi:10.1200/JCO.2016.70.3934
- 456 29. Diéras V, Miles D, Verma S, et al. Trastuzumab emtansine versus capecitabine plus  
457 lapatinib in patients with previously treated HER2-positive advanced breast cancer (EMILIA):  
458 a descriptive analysis of final overall survival results from a randomised, open-label, phase 3  
459 trial. *Lancet Oncol*. 2017;18(6):732-742. doi:10.1016/S1470-2045(17)30312-1
- 460 30. Meehl PE. *Clinical versus Statistical Prediction: A Theoretical Analysis and a Review*  
461 *of the Evidence*. Minneapolis, MN, US: University of Minnesota Press; 1954.  
462 doi:10.1037/11281-000
- 463 31. Dawes RM, Faust D, Meehl PE. Clinical Versus Actuarial Judgment. *Science*.  
464 1989;243:7. doi:10.1126/science.2648573
- 465 32. Grove WM, Zald DH, Lebow BS, Snitz BE, Nelson C. Clinical versus mechanical  
466 prediction: a meta-analysis. *Psychol Assess*. 2000;12(1):19-30.

- 467 33. Schroen AT, Petroni GR, Wang H, et al. Challenges to accrual predictions to phase III  
468 cancer clinical trials: a survey of study chairs and lead statisticians of 248 NCI-sponsored  
469 trials. *Clin Trials Lond Engl*. 2011;8(5):591-600. doi:10.1177/1740774511419683
- 470 34. London JW, Balestrucci L, Chatterjee D, Zhan T. Design-phase prediction of potential  
471 cancer clinical trial accrual success using a research data mart. *J Am Med Inform Assoc*  
472 *JAMIA*. 2013;20(e2):e260-266. doi:10.1136/amiajnl-2013-001846
- 473 35. e-CFR. *Electronic Code of Federal Regulations*. Vol Title 21 → Chapter I →  
474 Subchapter D → Part 312.; 1987. [https://www.ecfr.gov/cgi-bin/text-  
475 idx?SID=e05abc4c40756e7a19bd0ca2bc38b112&mc=true&node=pt21.5.312&rgn=div5](https://www.ecfr.gov/cgi-bin/text-idx?SID=e05abc4c40756e7a19bd0ca2bc38b112&mc=true&node=pt21.5.312&rgn=div5).  
476 Accessed February 12, 2019.
- 477 36. DiMasi JA, Reichert JM, Feldman L, Malins A. Clinical Approval Success Rates for  
478 Investigational Cancer Drugs. *Clin Pharmacol Ther*. 2013;94(3):329-335.  
479 doi:10.1038/clpt.2013.117
- 480 37. Smietana K, Siatkowski M, Møller M. Trends in clinical success rates. *Nat Rev Drug*  
481 *Discov*. 2016;15:379. doi:10.1038/nrd.2016.85
- 482 38. DiMasi J, Hermann J, Twyman K, et al. A Tool for Predicting Regulatory Approval  
483 After Phase II Testing of New Oncology Compounds. *Clin Pharmacol Ther*. 2015;98(5):506-  
484 513. doi:10.1002/cpt.194
- 485 39. Dahlin E, Nelson GM, Haynes M, Sargeant F. Success rates for product development  
486 strategies in new drug development. *J Clin Pharm Ther*. 2016;41(2):198-202.  
487 doi:10.1111/jcpt.12362
- 488 40. Schachter AD, Ramoni MF. Clinical forecasting in drug development. *Nat Rev Drug*  
489 *Discov*. 2007;6(2):107-108. doi:10.1038/nrd2246
- 490
- 491

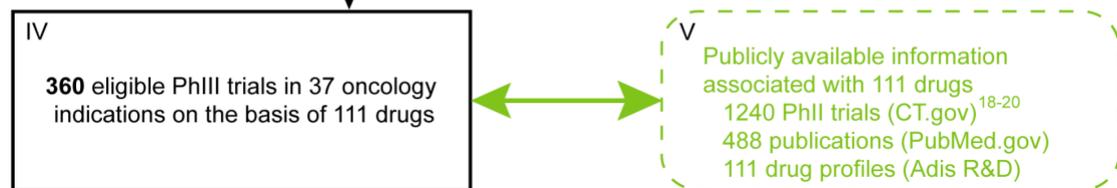
492 **Figures**

493 **Figure 1**

A) Screen



B) Data base creation



494

495

496 Consort flow diagram. (A) Screen. 360 eligible PhIII studies were identified screening

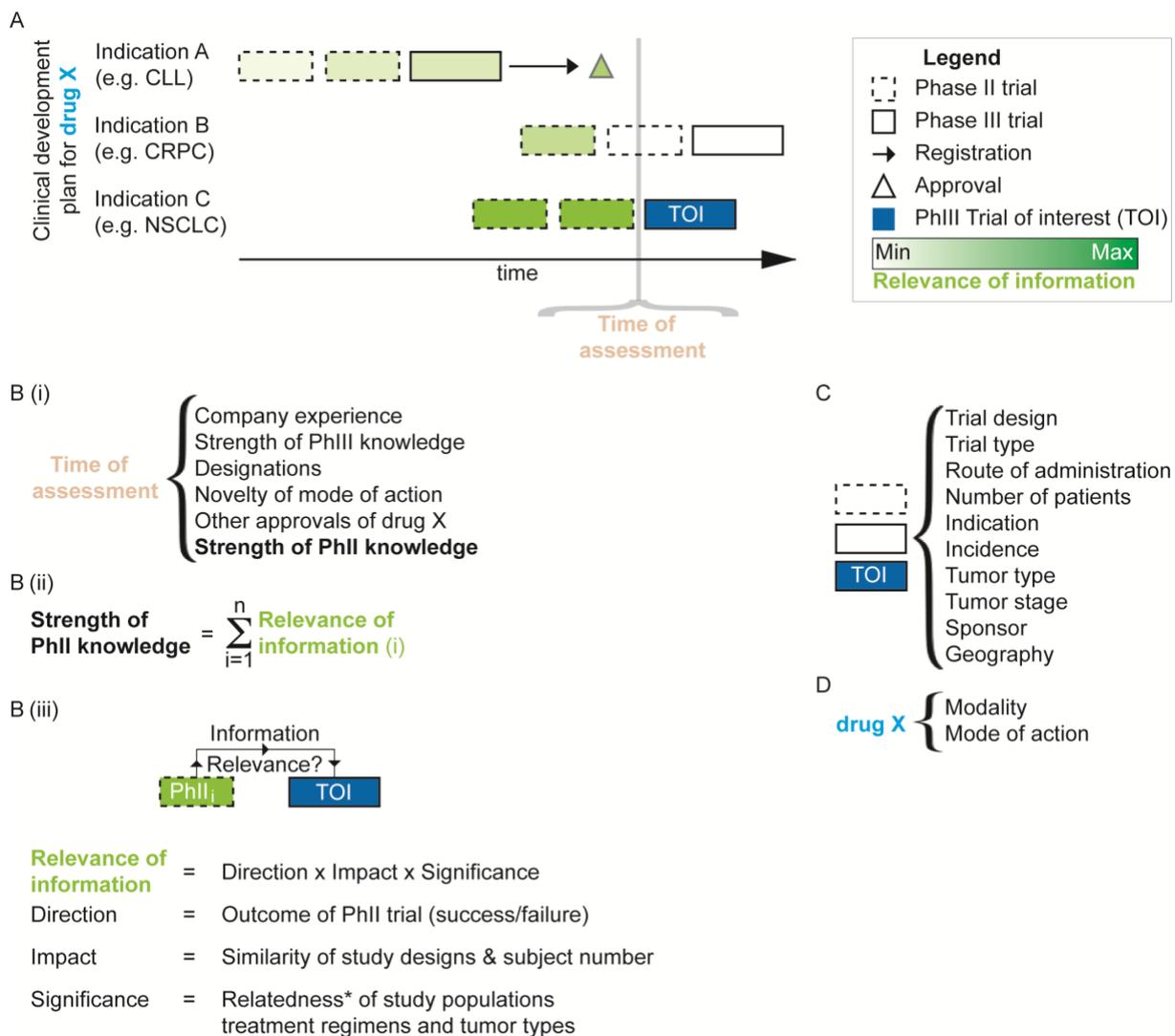
497 ClinicalTrials.gov<sup>18</sup> and other clinical trial registries<sup>19,20</sup> (I-IV, black font). (B) Database

498 creation. Associated information for each eligible PhIII trial was searched for and added

499 using public sources (V, green font) to create a project specific database<sup>18-22</sup> (IV).

500

501 **Figure 2**



502

503 Identification, classification and weighting of variables relevant for PhIII trial assessment. (A)

504 Schematic clinical development plan for drug X developed by sponsor Y. At time of

505 assessment (vertical grey line) not all trials are completed (boxes with white fillings) hence

506 cannot contribute information for assessing the PoS of the trial of interest (TOI, blue). Trials

507 that were completed prior to the assessment carry information, but the relevance of this

508 information (shades of green) varies with regard to the specific characteristics of the TOI.

509 Generally speaking, the closer the patient population, the study design and the treatment

510 algorithm of a given trial with regard to the TOI, the higher its relevance. Note, the clinical

511 development plan (CDP) illustrates only elements, which are considered in the database. We

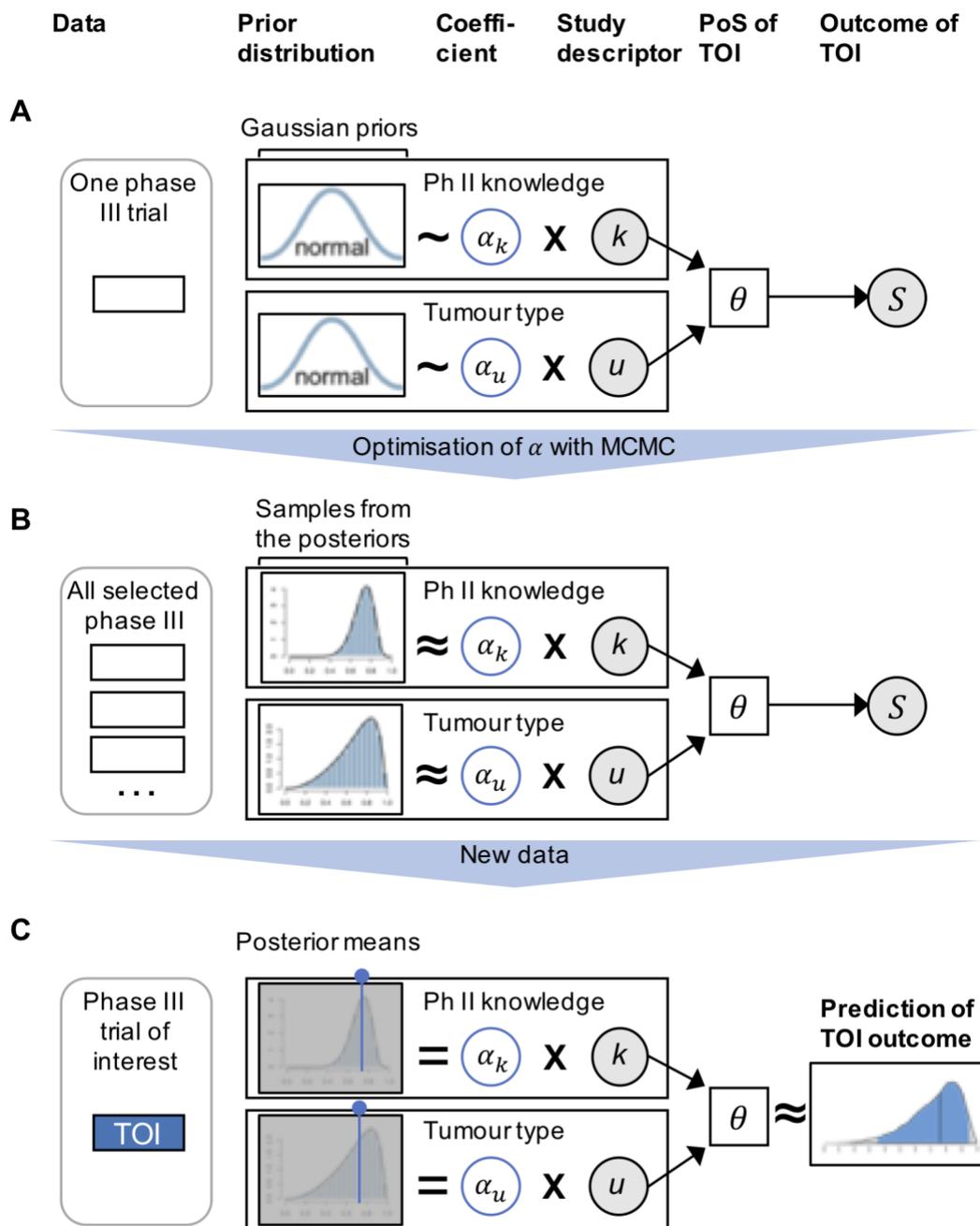
512 identified three large categories of parameters: Time-dependent variables (B), trial-specific

513 characteristics (C) and drug-related characteristics (D). (B.i) Time-dependent variables at

514 time of assessment. Each of these variables is a composite score of several sub-variables,  
515 as illustrated by the variable Strength of PhII knowledge (SPIIK, bold). (B.ii) Weighting of  
516 relevant information exemplified by variable 'Strength of phase II knowledge'. SPIIK is a  
517 composite of all relevant information (ROI) *i* carried by all PhII studies completed at least 2  
518 months prior to the time of assessment. The more PhII studies *i* were completed, the  
519 stronger the body of evidence. (B.iii) Definition of ROI. Any given PhII trial *i* (green)  
520 conducted with drug X (not shown) carries information (black arrows) that are potentially  
521 relevant for assessing the PoS of the TOI (blue). The ROI for a given trial *i* is defined as the  
522 product of three factors, each of which can be broken down into further subfactors, that may  
523 even be further broken down (e.g. Relatedness, \*see eFig 2 for details) until an objective and  
524 quantifiable level of information is found. Note, that there is a unique ROI for each  
525 combination of PhII and TOI, thus a unique SPIIK for each TOI. (C) Characteristics of novel  
526 therapeutic. (D) Inherent characteristics of trial of interest.

527

528 **Figure 3**



529

530 Schematic representation of Bayesian logistic regression. (A) Starting point: Bayesian model

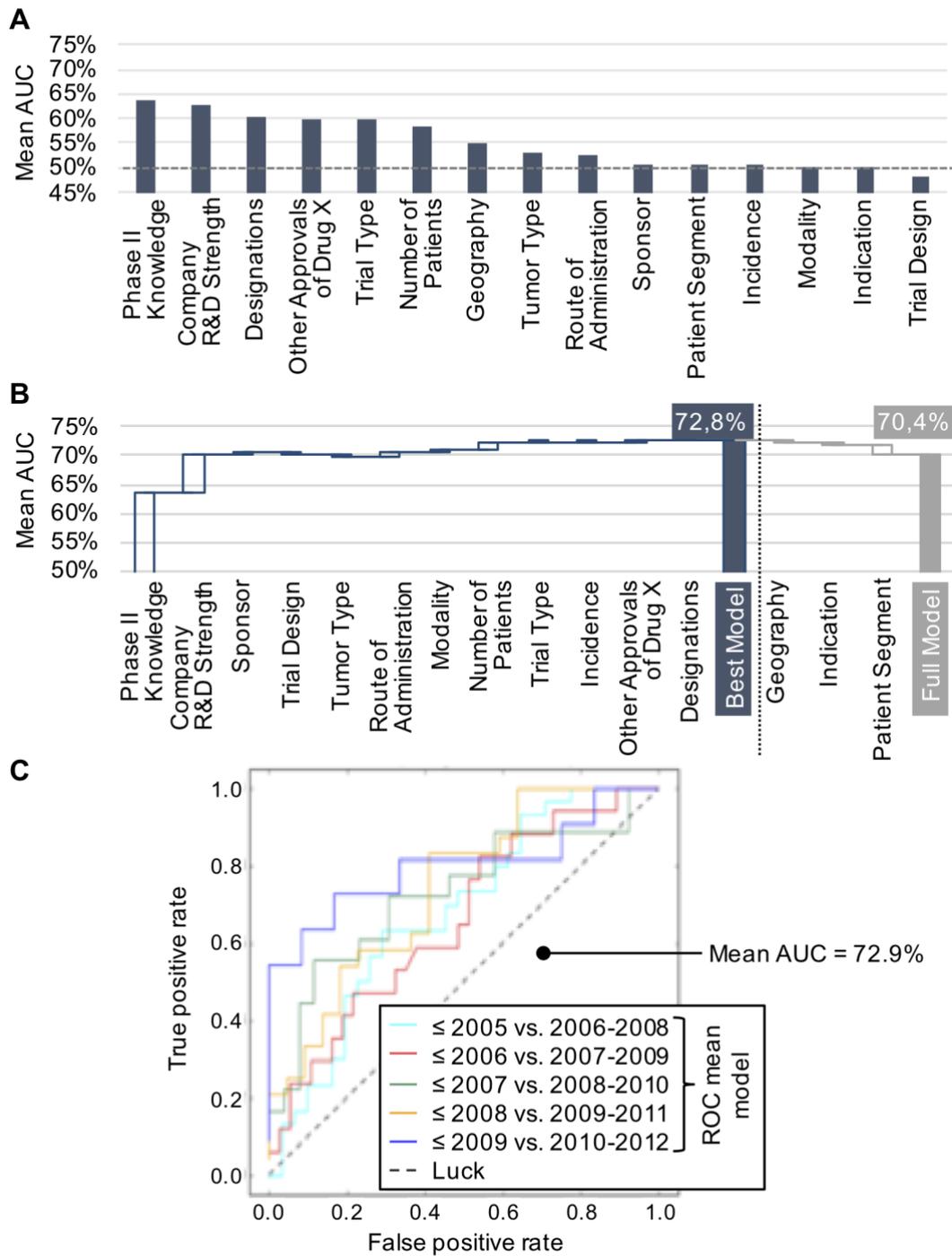
531 without assumptions about the parameters. (B) Model training: samples from the posteriors

532 gained from Markov chain Monte Carlo (MCMC). (C) Calculation of  $PoS_{TOI}$  for ongoing or

533 planned PhIII trial.

534

535 **Figure 4**



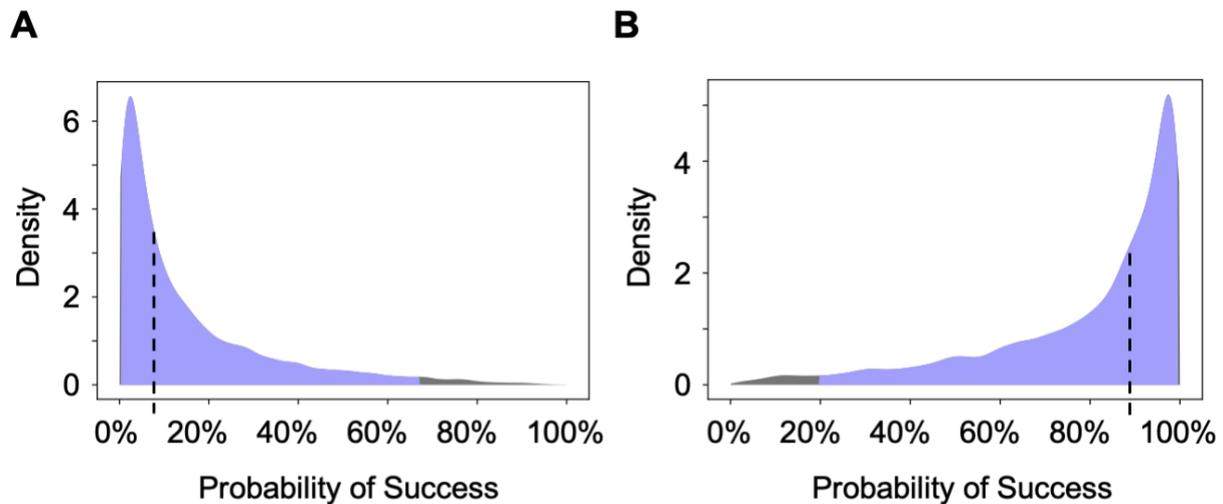
536

537 Model calibration and evaluation of performance. (A) AUC analysis of single-variables. The  
 538 greater the distance from AUC = 0.5 (dotted line), the stronger the predictive performance of  
 539 the variable. (B) At each step, the composite score performing best in combination with the  
 540 intermediate model is selected. The predictive performance is at its maximum after the  
 541 inclusion of 12 composite scores. The composite scores geography, indication, and patient

542 segment are not included in the model because they fail to increase its performance. (C) For  
543 the best model, the receiver operating characteristic curves are displayed for each time  
544 series. The overall model performance is given by the mean AUROC over all time points and  
545 is ~73%.

546

547 **Figure 5**



**C**

J Clin Oncol. 2017 Jul 10;35(20):2279-2287. doi: 10.1200/JCO.2016.70.3934. Epub 2017 May 24.

**Gefitinib and EGFR Gene Copy Number Aberrations in Esophageal Cancer.**

**D**

Lancet Oncol. 2017 Jun;18(6):732-742. doi: 10.1016/S1470-2045(17)30312-1. Epub 2017 May 16.

**Trastuzumab emtansine versus capecitabine plus lapatinib in patients with previously treated HER2-positive advanced breast cancer (EMILIA): a descriptive analysis of final overall survival results from a randomised, open-label, phase 3 trial.**

548

549 Selected PoS distributions. (A) The 13 composite scores selected for the model and the 58

550 variables making up those scores (unnamed for clarity) are displayed for two exemplary

551 clinical trials. The variables characterizing each trial are marked by a dot. The variables'

552 weight is color coded: hue conveys the sign and intensity the amplitude. Blue indicates a

553 positive influence on PoS prediction, red a negative one, and white no influence. (B) Study 1:

554 The predicted PoS of the phase III trial in esophageal cancer with gefitinib (NCT01243398) is

555 7.8%<sub>PoS</sub> (dotted line) and the 95%<sub>CI</sub> credible interval ranges from 0.3%<sub>PoS</sub> to 67.2%<sub>PoS</sub>

556 (colored area under the curve). (C) Study 2: The predicted PoS of the phase III trial in breast

557 cancer with trastuzumab emtansine (NCT00829166) is 88.6%<sub>PoS</sub> and the 95%<sub>CI</sub> credible

558 interval ranges from 19.6%<sub>PoS</sub> to 99.6%<sub>PoS</sub>. The predicted PoS for those two studies are

559 consistent with results published in peer-reviewed articles: the study with gefitinib failed (D),

560 while the study with trastuzumab emtansine was successful (E).