

1 **Title:** A Comparison of the Randomized Clinical Trial Efficacy and Real-World Effectiveness of Tofacitinib
2 for the Treatment of Inflammatory Bowel Disease: A Cohort Study

3

4 **Authors:** Vivek A. Rudrapatna MD^{1,2}, Benjamin S. Glicksberg PhD², Atul J. Butte MD^{2,3,4}

5

6 **Affiliations:**

7 1. Division of Gastroenterology, Department of Medicine, University of California, San Francisco, CA

8 2. Bakar Computational Health Sciences Institute, University of California, 550 16th Street, San

9 Francisco, CA

10 3: Department of Pediatrics, University of California, San Francisco, CA

11 4: Center for Data-Driven Insights and Innovation, University of California Health, Oakland, CA

12

13 **Corresponding Author:** Atul J. Butte. 550 16th Street, 4th Floor Box 0110, San Francisco, CA 94158-2549.

14 tel: 415-514-0511. email: atul.butte@ucsf.edu

15

16 **Abstract**

17 **Background:** Real-world data are receiving attention from regulators, biopharmaceuticals and payors as
18 a potential source of clinical evidence. However, the suitability of these data to produce evidence
19 commensurate with randomized controlled trials (RCTs) and the best practices in their use remain
20 unclear. We sought to compare the real-world effectiveness of Tofacitinib in the treatment of IBD
21 against efficacy rates published by corresponding RCTs.

22
23 **Methods:** Electronic health records at the University of California, San Francisco (UCSF) were queried
24 and reviewed to identify 86 Tofacitinib-treated IBD patients through 4/2019. The primary endpoint was
25 treatment effectiveness. This was measured by time-to-treatment-discontinuation and by the primary
26 endpoints of RCTs in Ulcerative Colitis (UC) and Crohn's Disease (CD). Endpoints were measured and
27 analyzed following a previously published protocol and analysis plan.

28
29 **Findings:** 86 patients (68 with UC, 18 with CD) initiated Tofacitinib for IBD treatment. Most of the data
30 needed to calculate baseline and follow-up disease activity indices were documented within the EHR
31 (77% for UC, 91% for CD). Baseline characteristics of the UCSF and RCT cohorts were similar, except for a
32 longer disease duration and 100% treatment failure of Tumor Necrosis Factor inhibitors in the former.
33 None of the UCSF cohort would have met the RCT eligibility criteria due to multiple reasons.

34
35 The rate of achieving the RCT primary endpoints were highly similar to the published rates for both UC
36 (16%, P=0.5) and CD (38%, P=0.8). However, treatment persistence was substantially higher: 69% for UC
37 (week 52) and 75% for CD (week 26).

38
39 **Interpretation:** An analysis of routinely collected clinical data can reproduce published Tofacitinib
40 efficacy rates, but also indicates far greater treatment durability than suggested by RCTs including
41 possible benefit in CD. These results underscore the value of real-world studies to complement RCTs.

42
43 **Funding:** The National Institutes of Health and UCSF Bakar Institute

44 **Research in Context**

45

46 **Evidence before this study**

47 Tofacitinib is the most recently approved treatment for Ulcerative Colitis. Data related to treatment
48 efficacy for either IBD subtype is generally limited, whether from controlled trials or real-world studies.
49 A search of clinicaltrials.gov was performed in January 2019 for completed phase 2 or 3, interventional,
50 placebo-controlled clinical trials matching the terms “Crohn’s Disease” OR “Ulcerative Colitis” in the
51 conditions field, and matching “Placebo” AND “Tofacitinib” OR “CP-690,550”) in the Interventions field.
52 We identified three Phase 3 trials for UC (OCTAVE trials, all initially reported in a single article in 2016)
53 and three Phase 2 trials of CD (two published in the same article in 2017, one reported in 2014). The
54 Phase 3 UC trials reported 57.6% pooled clinical response rate in the Tofacitinib-assigned groups after 8
55 weeks (induction), and a 37.5% pooled remission rate among eligible induction trial responders in the
56 Tofacitinib-assigned groups at 52 weeks. The 2017 CD trial reported a 70.8% pooled rate of response or
57 remission in the Tofacitinib-assigned groups after 8 weeks, and a 47.6% pooled rate of response or
58 remission among enrolled induction-trial responders at 26 weeks. A bias assessment of both UC and CD
59 trials indicated a high risk of attrition bias and unclear risk of bias related to conflicts of interest. We also
60 performed a search of pubmed.gov in January 2019 using search terms (“Colitis” OR “Crohn’s”) AND
61 (“Tofacitinib” OR “CP-690,550”) OR “real-world” to identify cohort studies of Tofacitinib efficacy in
62 routine clinical practice. No studies meeting these criteria were identified.

63

64 **Added value of this study**

65 This is one of the early studies to closely compare the results of clinical trials with the continuously-
66 updated data captured in the electronic health records, and the very the first to assess the efficacy-
67 effectiveness gap for Tofacitinib. We found that none of the patients treated at our center thus far
68 would have qualified for the clinical trial based on published eligibility criteria. We found that the drug
69 appeared to perform similarly to its efficacy when using the endpoints reported in clinical trials, but
70 treatment persistence was significantly greater than would have been expected from the reported trial
71 outcomes: 69% for UC at week 52 and 75% for CD at week 26.

72

73 **Implications of all the available evidence**

74 Tofacitinib is an effective treatment for the Ulcerative colitis and may be efficacious for Crohn’s disease.
75 Controlled trials may not be representative of real-world cohorts, may not be optimally designed to
76 identify efficacious drugs, and may not accurately predict patterns of use in clinical practice. Further
77 studies using real-world data as well as methods to enable their proper use are needed to confirm and
78 continuously monitor the efficacy and safety of drugs, both for on- and off-label use.

79 **Introduction:**

80

81 Inflammatory Bowel Disease (IBD) has increasingly been recognized as a global disease with accelerating
82 incidence and prevalence in newly industrialized nations.¹ Although IBD has historically been quite
83 morbid – associated with progressive bowel damage, malnutrition, chronic pain and neoplasia – recent
84 decades have seen a revolution in disease treatment and natural history with the advent of molecularly
85 targeted therapies.

86

87 Most of these new treatments have been monoclonal antibody biologics: agents that are expensive to
88 prepare and deliver, and associated with a loss of response over time. The first oral small-molecule for
89 the treatment of moderate to severe Ulcerative Colitis (UC) was approved by the EMA and the US FDA in
90 2018. Tofacitinib was approved on the basis of positive Phase 3 studies where it showed a statistically
91 significant reduction in the Mayo score over placebo by week 52.² Tofacitinib was also studied in two
92 Phase IIb randomized controlled trials (RCTs) of Crohn’s Disease (CD), the other major subtype of IBD. In
93 CD however, no statistical difference from placebo was seen in Crohn’s Disease Activity Index (CDAI)
94 reduction at week 26.³ Consequently, further investigation of Tofacitinib in Crohn’s Disease was
95 abandoned.

96

97 RCTs have long been considered the gold standard for clinical evidence.⁴ However, these studies have
98 come under greater scrutiny due to questions of their generalizability to the average clinical setting. In
99 particular, they have been associated with restrictive eligibility criteria that would exclude many real-
100 world patients being considered for care.⁵ Moreover, RCTs often do not measure the very same
101 endpoints as those used in routine clinical care. In practice the decision to continue, modify, or
102 discontinue treatment typically requires a patient-centric and setting-specific discussion of risks,
103 benefits, and alternatives. As such, RCT endpoints which are commonly binary, uniformly applied, and
104 analyzed according to the intention-to-treat principle are often less apt at answering common patient
105 questions such as: “What are the chances that this treatment is going to work?”

106

107 Recent years have seen significant interest in *real-world data* (RWD) in part due to their potential to
108 answer these sorts of pragmatic questions.⁶⁻⁸ However, their use has been fraught with many challenges
109 including missing data and misclassification. Widely-adopted standards for the handling of this data

110 have not been clearly established, nor is careful benchmarking against other gold-standard sources of
111 evidence presently a common practice.

112

113 In this cohort study we attempt to answer the following questions: Is routinely collected clinical data
114 complete enough to measure the endpoints assessed in IBD RCTs of Tofacitinib? Using both these
115 regulatory endpoints as well as the more pragmatic measure of *time to treatment discontinuation*, what
116 is the real-world effectiveness of Tofacitinib? How does it compare to trial efficacy rates? Are there any
117 systematic differences between the populations under study in both controlled and real-world settings?

118

119 **Methods:**

120

121 This study was performed in accordance with the STROBE statement⁹ (See Supplemental Content).

122

123 Patient identification and covariate extraction

124

125 We queried a structured database of deidentified Electronic Health Records (EHR) at the University of
126 California, San Francisco (UCSF) to identify patients meeting the following criteria: 1) age 18 or older, 2)
127 presence of a medication order for Tofacitinib, and 3) presence of an IBD diagnosis code (ICD-10-CM
128 K50*/K51*) assigned during a Gastroenterology clinic visit (Table 1). The scope of the search extended
129 from the instantiation of the EHR software (6/2012) through the time of the query (4/2019). We
130 obtained Institutional Review Board approval (#18-24588) to obtain identifiable patient record data, to
131 confirm that these patients had initiated treatment on Tofacitinib for the treatment of IBD, and to
132 assess treatment compliance. Compliance was defined as adherence to the treatment plan (e.g. dose,
133 frequency, duration) determined by the ordering provider.

134

135 Following a openly published protocol and statistical analysis plan,¹⁰ we reviewed all patient records to
136 measure the time to treatment discontinuation or last known use. We selected a subset of these records
137 to perform more detailed covariate extraction from the EHR following the aforementioned protocol.
138 These covariates included baseline demographics as well as the primary endpoints of the Phase 2b/3
139 RCTs of Tofacitinib for CD³ and UC² respectively. The time windows used for covariate extraction were
140 months -6 through 0 for the baseline data, and months 2 through 8 for the follow-up data. We selected
141 these windows in order to balance data availability and typical practice patterns with comparability to
142 protocol-driven endpoint measurements in RCTs.

143

144 Data quality and missing data assessments

145

146 We assessed the quality of the data in detail prior to proceeding with further analysis. We confirmed the
147 basic epidemiological properties of the dataset, such as the bimodal age of IBD onset which has been
148 previously reported^{11,12} (Supplemental Content eFigure 1). We annotated missing data and characterized
149 its distribution. Given that missing data was present across several variables following different non-
150 normal distributions, we performed multiple imputation by chained equations using random forest

151 classification/regression models (20 Markov chains, 10 iterations). We augmented our imputation model
152 with endpoints of interest in order to facilitate ‘congeniality’.¹³ We used all available variables as
153 auxiliary variables to satisfy the missing at random (MAR) assumption and increase imputation power.
154 We examined trace plots and strip plots to confirm Markov chain convergence and the plausibility of
155 imputed values respectively (Supplement).

156

157 Comparison to RCT endpoints

158

159 We abstracted RCT endpoint definitions and rates from Sandborn et al., 2017² and Panés et al., 2017³
160 (Supplemental Methods). Because both trials reported endpoint rates among those with a favorable
161 response to treatment induction, we recalculated the overall maintenance endpoint rate as the product
162 of the induction and maintenance response rates. All endpoints were calculated from the active
163 treatment arm without attenuation by the placebo rate.

164

165 Statistics and computing

166

167 We performed point estimation and hypothesis testing by calculating Wald test statistics with pooled
168 standard errors following “Rubin’s rules”.¹⁴ We estimated the *time to treatment discontinuation* survival
169 distributions using the product-limit estimator. No competing events (e.g. mortality) were observed.
170 Treatment discontinuation due to loss of insurance coverage, as well as relocation or other lost-to-
171 follow-up events were rare and were treated as non-informatively censored events.

172

173 We performed statistical computing in the *R* statistical computing environment (3.6.0) using the
174 following packages: *pacman*, *data.table*, *survival*, *survminer*, *readxl*, *tidyr*, *scales*, *binom*, *ggplot2*,
175 *lubridate*, *randomForest*, *RMarkdown*, *mice*, *Hmisc*, and *magrittr* (Supplemental References). The
176 statistical code was independently reviewed by a co-author (BSG). Synthetic data and analysis files were
177 version-controlled using *Docker*.

178

179 Role of the funding source

180 The funding source had no role in any aspect of this study, including design, collection, analysis, data
181 interpretation, writing, nor the decision to submit this manuscript for publication.

182

183

184 **Results:**

185

186 Cohort Identification

187

188 We queried an EHR structured database system comprising ~1.2 million patient records to identify adult
189 patients with an IBD diagnosis assigned by a gastroenterologist and a medication order for Tofacitinib.
190 115 potential patient records were identified. We performed manual review to confirm that 86 patients
191 – 68 with Ulcerative Colitis (UC) and 18 with Crohn’s Disease (CD) – had initiated Tofacitinib specifically
192 to treat IBD (Figure 1). The other 29 patients were excluded during this process for multiple reasons,
193 including failure to start treatment due to payor denial, the decision to forgo the ordered medical
194 treatment in favor of surgery, and treatment initiated by a non-gastroenterologist for another
195 autoimmune condition. Non-compliance with Tofacitinib was rare (4%) in this cohort.

196

197 Data completeness assessment

198

199 We first sampled a subset of these 86 records to quantify the completeness of routinely collected
200 clinical data in the capture of IBD-relevant clinical indices. We identified 87% and 91% of all the UC and
201 CD data elements (demographics, indices) as available within the EHR. Within the set of data elements
202 needed to calculate the baseline and follow-up Mayo Score and CDAI (for UC and CD respectively), we
203 observed 77% and 91% data completeness. We performed multiple imputation on these missing
204 elements to enable further downstream analysis. For instance, patients deemed treatment failures on
205 the basis of endoscopic worsening, physician global assessment, and stool frequency but missing a rectal
206 bleeding subscore required imputation in order to calculate the full Mayo score.

207

208 Baseline cohort comparison

209

210 The baseline demographics of the subjects under study in the RCTs and the UCSF cohort were largely
211 similar (Table 1). Notable differences include the universal failure of TNF inhibitors in the UCSF cohort,
212 as well as a longer duration of disease in the UC patients. Patient groups had similar baseline Mayo
213 scores, c-reactive protein levels, and rates of corticosteroid use.

214

215 We assessed the proportion of UC patients who would have satisfied the eligibility criteria of the
216 corresponding Phase III RCT². We found that 0% of the UC patients initiated on Tofacitinib met these
217 criteria. The reasons for this were multifactorial (Table 2) but include prior use of Vedolizumab within
218 the previous year, use of high-dose prednisone or intravenous methylprednisolone at the time of
219 treatment initiation, possibility of requiring surgery during the treatment period, and the lack of
220 protocolized screening or scheduled steroid tapering in the routine clinical setting.

221
222 We separately explored what proportion of patients met the specific RCT entry criteria defined by the
223 Mayo score and CDAI for UC and CD respectively. 93% (73-98) of the UC patients had an eligible baseline
224 Mayo score, whereas 50% (19-82) of the CD patients had a baseline CDAI within the eligible range of the
225 corresponding RCT.

226 227 Efficacy vs Effectiveness

228
229 Time to treatment discontinuation analysis on the full cohort revealed nearly identical survival
230 distributions irrespective of IBD disease subtype (Figure 2). The overall probability of incident users
231 maintaining long-term treatment on Tofacitinib was 68% (58-80%). All failure events occurred within the
232 first seven months; among continued responders by month six, the probability of sustained long-term
233 response was 94%. Of note, the first use of the Tofacitinib occurred in 2013, and the longest duration of
234 effectiveness data relevant to treatment maintenance was 3.7 years.

235
236 22% of all subjects participating in the induction phase of the UC RCT² met the primary maintenance
237 endpoint of week 52 clinical remission. We observed a similar rate (16%) in the corresponding UCSF
238 cohort (6-37%, p-value 0.5). Similarly, the rate of achieving the primary endpoint in the CD RCT³ (34%)
239 was essentially the same as the point estimate of the real-world cohort (38%, p-value 0.8).

240
241 However, the rates of meeting these regulatory endpoints were substantially different from the week 52
242 and week 26 survival probabilities for UC and CD respectively (p-values 3e-14 and 9e-5). This difference
243 from the empiric response rate remained even when compared to the least stringent secondary
244 endpoint for the UC RCT (week 52 clinical response), met by 33% of all induction participants (p-value
245 0.05).

246

247 **Discussion:**

248

249 RWD has been receiving growing interest from a variety of parties in recent years. The European
250 Medicines Agency¹⁵ and the US Food and Drug Administration⁷ have been formalizing regulatory
251 pathways to evaluate RWD in support of new drug indications and post-marketing surveillance. RWD is
252 also being used at earlier stages of clinical development by biopharmaceuticals in order to improve drug
253 development and clinical trial design. Healthcare payors are beginning to use RWD to support
254 outcomes-based pricing contracts. Providers and patients are increasingly interested in understanding
255 how RWD can enable personalized medicine in clinical practice.

256

257 Although these developments have been promising, multiple barriers have precluded the otherwise
258 widespread adoption of RWD to improve healthcare. Missingness is ubiquitous in clinical data and
259 especially so in the EHR, but is typically either deemphasized or handled using ad hoc and biased
260 methods.¹⁶ Studies of large real-world clinical datasets typically depend on administratively coded
261 structured data which are prone to mis-annotation.¹⁷ The use of registered protocols, published analysis
262 plans, independent code review, and the public release of reproducible analysis documents are
263 uncommon in observational studies, and have generally contributed to their decreased credibility in
264 comparison to their better-funded, controlled study counterparts. Lastly, real-world studies have
265 prioritized pragmatic endpoints but without sufficient attention to careful benchmarking against
266 controlled trials using clinical instruments, further contributing to the credibility gap. Overall, there is a
267 lack of widely-adopted standards for the proper handling of RWD to robustly support the generation of
268 new clinical knowledge.

269

270 This study attempts to address these shortcomings and critically assess the potential of RWD to inform
271 both research and practice via the use-case of IBD. To curtail the potential bias from unblinded record
272 review, we published a detailed protocol and statistical analysis plan on an open platform in advance of
273 this work. In addition to capturing the pragmatic endpoint of *time to treatment discontinuation*, we
274 explicitly prioritized the benchmarking of our data to the corresponding RCTs and took special efforts to
275 capture the very same clinical instruments and endpoints assessed by regulatory bodies. We critically
276 assessed the missingness of our data and addressed it using principled methods. We performed an
277 independent code review and have prepared a synthetic data set and dynamic analysis documents to
278 maximize both the reproducibility of our results and the reusability of our code in other health systems.

279

280 The results of our study have important implications for both the future of real-world evidence studies
281 and specifically the treatment of IBD. We show that routinely collected clinical data from the EHR is
282 sufficiently rich to enable the measurement of the primary clinical instruments used in IBD RCTs with
283 limited missingness. We also demonstrate that point-estimates of clinical effectiveness derived from
284 RWD align well with those suggested by RCTs when using the same clinical instruments. These results
285 support the potential for robust inference from RWD.

286

287 However, our data also speak to the presence of a significant efficacy-effectiveness gap between
288 treatment durability in practice and regulatory endpoints measured in clinical trials. More than the
289 regulatory endpoint, these pragmatic endpoints may provide the best answer to patients who ask: “How
290 likely is this drug to work for me?” (Answer: “Over half of patients like you will respond to treatment
291 long-term. It can take six months on treatment to know if you fall in this category”). Similarly, this
292 pragmatic endpoint may be most relevant to payors increasingly keen on paying for value and
293 anticipating the long-term costs of care.

294

295 Additionally, our study speaks to a number of implications on current trial design in IBD and their
296 application to practice. We found that only half of the CD patients at our center mounted a CDAI score
297 within range of the corresponding RCT. The CDAI score is known to correlate poorly with objective
298 markers of mucosal inflammation,^{18,19} and these results overall suggest the need for improved clinical
299 instruments for CD.²⁰ We also found that 100% of UC patients at our center were “functionally
300 disqualified” from the corresponding RCTs, consistent with prior work.⁵ On our review most of the
301 reasons for this were related to the careful control of alternative explanations to any differences in trial
302 endpoints seen between treatment arms. These observations reflect the fundamental trade-off
303 between internal and external validity, and specifically the importance of integrating controlled- and
304 real-world studies to obtain the highest quality clinical evidence.

305

306 Perhaps the most intriguing result of this study is the finding of identical survival distributions for the
307 real-world use of Tofacitinib in CD and UC. Given the positive trial results of Tofacitinib for UC and two
308 negative trial results in CD, we expected to see a divergence of survival distributions and a clear signal
309 that the drug does not work in CD. This finding to the contrary suggests two important corollaries: 1) the
310 CDAI may not be the optimal instrument to measure Crohn’s Disease activity, and 2) that the current

311 taxonomy of IBD (UC vs CD) may not be well rooted in molecular pathogenesis. Indeed, nearly all
312 pharmacological treatments between CD and UC are shared. Our results are consistent with recent post-
313 hoc analyses of the CD clinical trials suggesting a signal for efficacy if analyzed using more objective
314 endpoints,²¹ as well as the continued investigation of other JAK inhibitors for CD. These data underscore
315 the potential of real-world studies to harness “the wisdom of the masses,” to identify robust efficacy
316 signals that may support label expansion, and to expand options for the truly refractory patient facing
317 limited alternatives.

318
319 We acknowledge several limitations to this study. First, our sample size is still small. This is due to the
320 relatively recent approval of the drug for UC and its lack of EMA or FDA approval for CD. Although small
321 samples decrease the power to detect differences, the signals seen here were strong enough to reveal
322 many important effects, including survival distributions that are highly consistent with prior reports.²² Of
323 course, patient cohorts will continue to expand as this drug continues to be used, and real-world data
324 can keep growing.

325
326 Second, although we attempted to comprehensively characterize possible biases and place our findings
327 in context (Table 3), we cannot exclude the possibility of residual bias. Third, the validity of imputation
328 depends on the missing at random (MAR) assumption. We would argue that this assumption is plausible:
329 the most common reasons for missingness were clear indications of treatment failure (e.g. no rationale
330 to pursue endoscopy), there was complete measurement of the survival outcome, and the list of
331 auxiliary variables was extensive. Nevertheless, it remains fundamentally untestable. Lastly, our chart
332 review process – albeit guided by a detailed pre-published protocol – may not be entirely objective nor
333 scalable.

334
335 In summary, our study suggests that RWD can be a robust source of valuable clinical knowledge that
336 complements evidence from controlled trials. Tofacitinib appears to work as well as has been suggested
337 by clinical trials when using the same clinical instruments, but has greater real-world durability. Despite
338 negative trial results, Tofacitinib and JAK inhibitors more generally may be valuable for the treatment of
339 CD. Lastly, because clinical trials functionally exclude many patients, real-world studies are indispensable
340 to ensure the generalizability of RCT findings and add to the best evidence for clinical care.

341

342 **Author Contributions:** VAR and AJB conceived the project. VAR designed the chart review protocol,
343 performed chart extraction, conducted statistical analyses, and drafted this manuscript. BSG performed
344 code review and critically edited this manuscript. AJB supervised the project and critically edited this
345 manuscript.

346
347 **Acknowledgements:** The authors thank the UCSF Academic Research Services and Clinical Data
348 Research Consultation services for clinical informatics support. The authors would like to acknowledge
349 Dana Ludwig for his help in deidentifying and interpreting the UCSF EHR.

350
351 **Data Sharing Plan:** The analytic code in the form of a R markdown file as well as the accompanying data
352 set needed to reproduce the analysis in this work are available in a *Docker* container and will be
353 released for public use on <https://datadryad.org> at the time of publication to all investigators without
354 restriction (<https://doi.org/10.7272/Q6PZ5715>). These individual participant data were de-identified to
355 comply with the US Department of Health and Human Services 'Safe Harbor' guidance and applicable
356 laws and regulations concerning privacy and/or security of personal information. The data dictionary is
357 documented within the study protocol section of Supplemental Content.

358
359 **Financial Support:** Research reported in this publication was supported by funding from the UCSF Bakar
360 Computational Health Sciences Institute and the National Center for Advancing Translational Sciences of
361 the National Institutes of Health under award number UL1 TR001872. VAR was supported by the
362 National Institute of Diabetes and Digestive and Kidney Disease of the National Institutes of Health grant
363 under award number T32 DK007007-42. Its contents are solely the responsibility of the authors and do
364 not necessarily represent the official views of the NIH.

365
366 **Declaration of Interests:** No conflicts relevant to this publication exist.

367

368 **Figure Captions:**

369

370 **Figure 1:** Cohort Selection Schematic

371 **Figure 2:** Time to Tofacitinib Discontinuation. Shaded regions correspond to the 95% confidence
372 interval.

373

374 **Table Captions:**

375

376 **Table 1:** Demographic comparison of subjects studied in RCTs of Ulcerative Colitis and Crohn's Disease
377 with subjects assigned to receive Tofacitinib and corresponding patients at UCSF. *: *Includes patients*
378 *who received intravenous corticosteroids at the time of Tofacitinib initiation.*

379

380 **Table 2:** Most common reasons disqualifying the real-world ulcerative colitis cohort from meeting the
381 OCTAVE trial eligibility criteria

382

383 **Table 3:** Qualitative Analysis of Bias

384

385 **References:**

- 386 1. Ng SC, Shi HY, Hamidi N, et al. Worldwide incidence and prevalence of inflammatory bowel
387 disease in the 21st century: a systematic review of population-based studies. *Lancet (London,*
388 *England)*. 2018;390(10114):2769-2778. doi:10.1016/S0140-6736(17)32448-0
- 389 2. Sandborn WJ, Su C, Sands BE, et al. Tofacitinib as Induction and Maintenance Therapy for
390 Ulcerative Colitis. *N Engl J Med*. 2017;376(18):1723-1736. doi:10.1056/NEJMoa1606910
- 391 3. Panés J, Sandborn WJ, Schreiber S, et al. Tofacitinib for induction and maintenance therapy of
392 Crohn's disease: results of two phase IIb randomised placebo-controlled trials. *Gut*.
393 2017;66(6):1049-1059. doi:10.1136/gutjnl-2016-312735
- 394 4. Straus SE, Glasziou P, Richardson WS, Haynes RB. *Evidence-Based Medicine: How to Practice and*
395 *Teach EBM.*; 2018.
- 396 5. Ha C, Ullman TA, Siegel CA, Kornbluth A. Patients Enrolled in Randomized Controlled Trials Do
397 Not Represent the Inflammatory Bowel Disease Patient Population. *Clin Gastroenterol Hepatol*.
398 2012;10(9):1002-1007. doi:10.1016/j.cgh.2012.02.004
- 399 6. Medicines Agency E. *An Agency of the European Union Regulatory Perspective on Real World*
400 *Evidence (RWE) in Scientific Advice EMA Human Scientific Committees' Working Parties with*
401 *Patients' and Consumers' Organisations (PCWP) and Healthcare Professionals' Organisations*
402 *(HCPWP)*.
- 403 7. *Framework for FDA's Real-World Evidence Program.*; 2018.
404 <https://www.fda.gov/media/120060/download>. Accessed June 26, 2019.
- 405 8. Real-world evidence: From activity to impact in healthcare decision making | McKinsey.
406 [https://www.mckinsey.com/industries/pharmaceuticals-and-medical-products/our-insights/real-](https://www.mckinsey.com/industries/pharmaceuticals-and-medical-products/our-insights/real-world-evidence-from-activity-to-impact-in-healthcare-decision-making)
407 [world-evidence-from-activity-to-impact-in-healthcare-decision-making](https://www.mckinsey.com/industries/pharmaceuticals-and-medical-products/our-insights/real-world-evidence-from-activity-to-impact-in-healthcare-decision-making). Accessed June 13, 2019.
- 408 9. von Elm E, Altman DG, Egger M, et al. The Strengthening the Reporting of Observational Studies
409 in Epidemiology (STROBE) Statement: guidelines for reporting observational studies. *Int J Surg*.
410 2014;12(12):1495-1499. doi:10.1016/j.ijsu.2014.07.013
- 411 10. Rudrapatna VA, Butte AJ. Robust measurement of the real world effectiveness of Tofacitinib for
412 the treatment of Ulcerative Colitis using electronic health records: a protocol and statistical
413 analysis plan. doi:10.17504/protocols.io.2bqgamw
- 414 11. Ekblom A, Helmick C, Zack M, Adami HO. The epidemiology of inflammatory bowel disease: a
415 large, population-based study in Sweden. *Gastroenterology*. 1991;100(2):350-358.
416 <http://www.ncbi.nlm.nih.gov/pubmed/1985033>. Accessed June 20, 2019.

- 417 12. Bernstein CN, Wajda A, Svenson LW, et al. The epidemiology of inflammatory bowel disease in
418 Canada: a population-based study. *Am J Gastroenterol*. 2006;101(7):1559-1568.
419 doi:10.1111/j.1572-0241.2006.00603.x
- 420 13. Meng X-L. Multiple-Imputation Inferences with Uncongenial Sources of Input. *Stat Sci*.
421 1994;9(4):538-558. doi:10.1214/ss/1177010269
- 422 14. Rubin DB, Wiley J, York N, Brisbane C, Singapore T. *Multiple Imputation for Nonresponse in*
423 *Surveys*; 1987. <https://www.onlinelibrary.wiley.com/doi/pdf/10.1002/9780470316696.fmatter>.
424 Accessed June 20, 2019.
- 425 15. Cave A, Cerreta F. *Use of Real World Data in Development Programmes*.; 2017.
426 [https://www.ema.europa.eu/en/documents/presentation/presentation-use-real-world-data-](https://www.ema.europa.eu/en/documents/presentation/presentation-use-real-world-data-development-programmes-dr-alison-cave-dr-francesca-cerreta_en.pdf)
427 [development-programmes-dr-alison-cave-dr-francesca-cerreta_en.pdf](https://www.ema.europa.eu/en/documents/presentation/presentation-use-real-world-data-development-programmes-dr-alison-cave-dr-francesca-cerreta_en.pdf). Accessed June 26, 2019.
- 428 16. Buuren S van. *Flexible Imputation of Missing Data*.
- 429 17. Rudrapatna VA, Glicksberg BS, Avila P, Harding-Theobald E, Wang C, Butte AJ. Accuracy of
430 Medical Billing Data Against the Electronic Health Record in the Measurement of Colorectal
431 Cancer Screening Rates. *medRxiv*. January 2019:19004598. doi:10.1101/19004598
- 432 18. Jones J, Loftus E V., Panaccione R, et al. Relationships Between Disease Activity and Serum and
433 Fecal Biomarkers in Patients With Crohn's Disease. *Clin Gastroenterol Hepatol*. 2008;6(11):1218-
434 1224. doi:10.1016/J.CGH.2008.06.010
- 435 19. Ricanek P, Brackmann S, Perminow G, et al. Evaluation of disease activity in IBD at the time of
436 diagnosis by the use of clinical, biochemical, and fecal markers. *Scand J Gastroenterol*.
437 2011;46(9):1081-1091. doi:10.3109/00365521.2011.584897
- 438 20. Peyrin-Biroulet L, Panés J, Sandborn WJ, et al. Defining Disease Severity in Inflammatory Bowel
439 Diseases: Current and Future Directions. *Clin Gastroenterol Hepatol*. 2016;14(3):348-354.e17.
440 doi:10.1016/J.CGH.2015.06.001
- 441 21. Sands BE, Panés J, Higgins PDR, et al. 14 POST-HOC ANALYSIS OF TOFACITINIB CROHN'S DISEASE
442 PHASE 2 INDUCTION EFFICACY IN SUBGROUPS WITH BASELINE ENDOSCOPIC OR BIOMARKER
443 EVIDENCE OF INFLAMMATION. *Gastroenterology*. 2018;154(1):S81.
444 doi:10.1053/j.gastro.2017.11.203
- 445 22. Weisshof R, Aharoni Golan M, Sossenheimer PH, et al. Real-World Experience with Tofacitinib in
446 IBD at a Tertiary Center. *Dig Dis Sci*. 2019;64(7):1945-1951. doi:10.1007/s10620-019-05492-y
447

UCSF EHR Research Database

n=1.2M



≥1 IBD ICD code assigned by a GI provider
+ ≥1 Tofacitinib medication order

n=115



Prescribed by GI to treat IBD +
Confirmed start on treatment

n=86



Ulcerative Colitis

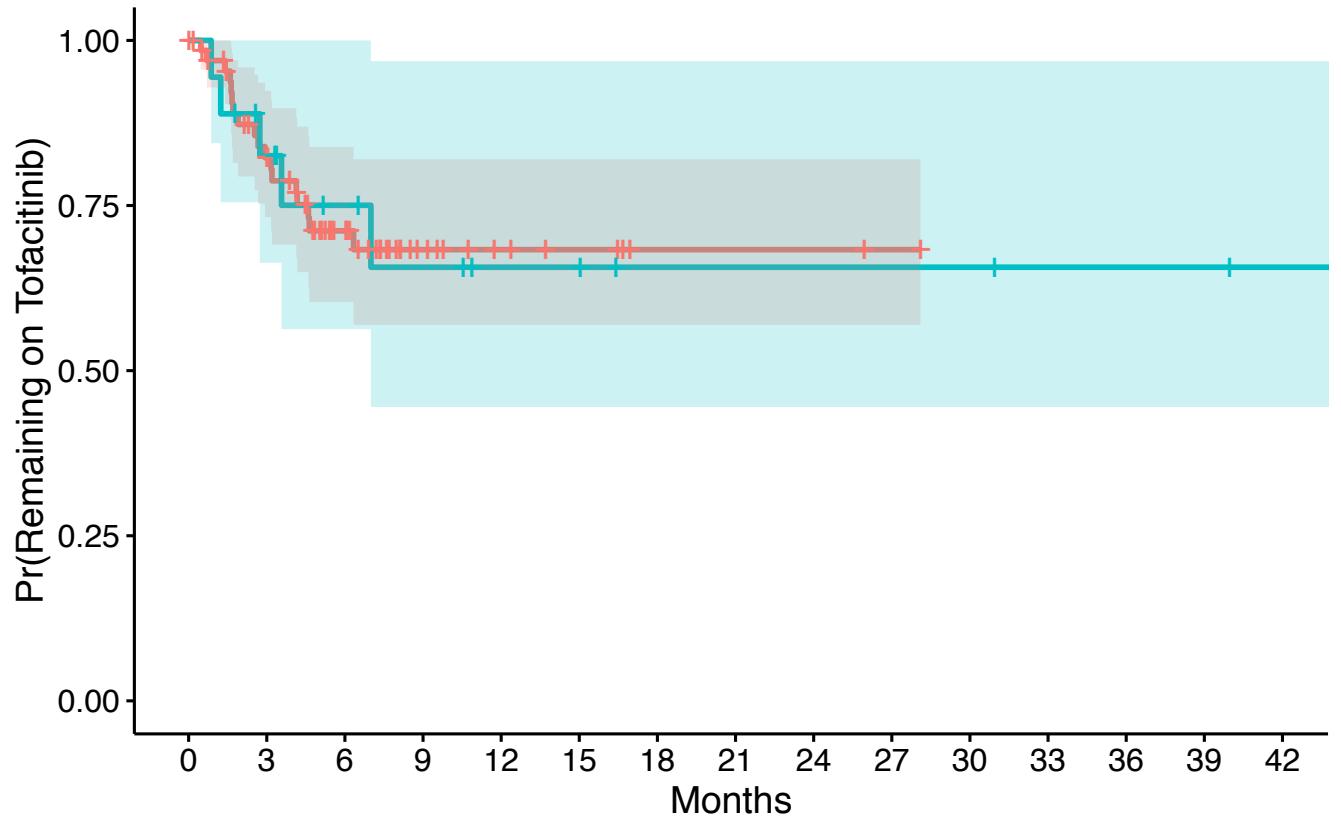
n=68

Crohn's Disease

n=18

Durability of Tofacitinib Treatment

Diagnosis UC CD



Number at risk

Diagnosis	0	3	6	9	12	15	18	21	24	27	30	33	36	39	42
UC	68	49	28	12	7	5	2	2	2	1	0	0	0	0	0
CD	18	13	9	7	5	5	3	3	3	3	3	2	2	2	1

	OCTAVE Induction 1 (n=475)	Sample of UC Cohort (n=28)		Panés et al. 2017 (n=86)	Sample of CD Cohort (n=13)
Male sex - no. (%)	277 (58.2)	16 (57)	Female, n (%)	47 (54.7)	9 (69.2)
Age – yr	41.3 ± 14.1	43.2 ± 14.4	Age, years		
Duration of disease - yr			-- Mean (SD)	39.3 (13.7)	39.7 (19.5)
-- Median	6.5	10.2	Weight, kg		
-- Range	0.3-42.5	2.2-51.4	-- Mean (SD)	71.6 (18.8)	69.9 (16.3)
Extent of Disease - no./total no. (%)			Race, n (%)		
-- Proctosigmoiditis	64/475 (13.7)	3/28 (10.7)	-- White	72 (83.7)	9 (69.2)
-- Left-sided colitis	158/475 (33.3)	6/28 (21.4)	-- Black	2 (2.3)	0 (0)
-- Extensive colitis or pancolitis	252/475 (53.1)	19/28 (67.9)	-- Asian	11 (12.8)	1 (7.7)
			-- Others	1 (1.2)	0 (0)
Total Mayo score	9.0 ± 1.4	8.5 ± 1.8	Duration since CD diagnosis, years		
Partial Mayo score	6.3 ± 1.2	6 ± 1.6	-- Mean (SD)	11.3 (9.7)	14.4 (8.2)
C-reactive protein - mg/liter			Extent of disease, n (%)		
-- Median	4.4	5.8	-- L1 (I/TI)	7 (8.1)	1 (7.7)
-- Range	0.1-208.4	0.8-70.6	-- L1/4 (I/TI + UGI)	2 (2.3)	2 (15.4)
Glucocorticoid use at baseline*	214 (45.0)	17 (60.7)	-- L2 (C)	5 (5.8)	0 (0)
Previous treatment with TNF antagonist – no. (%)	254 (53.4)	28 (100)	-- L2/4 (C + UGI)	16 (18.6)	1 (7.7)
Previous treatment failure - no. (%)			-- L3 (IC)	15 (17.4)	4 (30.8)
-- TNF antagonist	243 (51.1)	28 (100)	-- L3/4 (IC + UGI)	39 (45.3)	5 (38.5)
-- Glucocorticoid	350 (73.5)	24 (85.7)	Prior use of TNFi, n (%)	66 (76.7)	13 (100)
-- Immunosuppressant	360 (75.6)	21 (75)	Use of corticosteroids at study entry, n (%)	28 (32.6)	7 (53.8)
			Baseline CDAI score		
			-- Mean (SD)	320 (61.66)	374 (183.73)
			Baseline CRP, mg/L		
			-- Median (min-max)	5.5 (0.2-126)	28.7 (3.5-107)

Most common reasons disqualifying the real-world UC cohort from meeting the OCTAVE trial inclusion criteria:

Recent/Current use of other prohibited immunosuppressives:

- Vedolizumab use within the past 1 year
- Prednisone at doses exceeding 25mg/d (or equivalent)
- Anti-Tumor Necrosis Factor-Alpha use within the last 8 weeks
- Immunomodulator use within the past 2 weeks
- Any corticosteroid or 5-aminosalicylate per rectum within past 2 weeks
- Intravenous corticosteroid use within the past 2 weeks

Likely to require surgery during treatment period

Fluctuating dose of allowable concomitant medications (e.g. oral corticosteroid) during induction period

History of prior surgery potentially affecting drug absorption

Lack of mandatory screening-period exams

- Negative pregnancy testing and confirmation of highly effective contraceptive use
- Exclusion of colorectal dysplasia within past year
- Stool Culture

	Description	Effect on efficacy/effectiveness	Potential Solutions
Confounding Bias (e.g. Bias by Indication)			
<ul style="list-style-type: none"> Influences on decision to treat based on likelihood of response 	Many real-world patients initiated Tofacitinib following hospitalization, surgical consultation and intravenous steroids (all RCT exclusion criteria).	Possible decreased real-world effectiveness	Model-based approaches (e.g. propensity-score matching, inverse probability of treatment weighting)
Selection Bias			
<ul style="list-style-type: none"> Restrictive RCT eligibility criteria (See Table 2) 	Multiple exclusion criteria may limit generalizability to ordinary populations	Unclear effect on treatment efficacy measured by RCTs	Further studies of real-world cohorts
<ul style="list-style-type: none"> Tertiary care referral center population, off-label use in refractory patients 	Greater prevalence of treatment-resistant patients	Possible decreased real-world effectiveness	Further studies of real-world cohorts
Measurement Bias			
<ul style="list-style-type: none"> 'Laxity' in real-world use patterns compared to protocol-driven endpoints 	Lack of steroid-de-escalation protocols and refractory patients with limited alternatives may tolerate a suboptimal treatment response for a longer period of time.	Increased real-world effectiveness measured by time to treatment discontinuation	Trend real-world response rates over time
<ul style="list-style-type: none"> Intention to treat analysis with non-responder imputation 	Subjects non-adherent to the active arm and/or follow-up protocols are analyzed as treatment non-responders	Decreased treatment effect measured by RCTs	Re-analysis of clinical trial data with use of alternative imputation methods
<ul style="list-style-type: none"> Use of 6 month windows for covariate capture 	Baseline and follow-up covariates were captured by protocol to accommodate real-world practice	Unclear effect on real-world effectiveness	Model-based approaches to more accurately estimate values at fixed time points
<ul style="list-style-type: none"> Unblinded chart review 	Chart reviewer not blinded to outcome	Unclear effects on real-world effectiveness	Use of multiple blinded reviewers (assessing unrelated components) with measurement of interrater reliability