

3 VisHikers' Guide to Evaluation: Competing 4 Considerations in Study Design

5
6 Emily Wall , Emory University, Atlanta, GA, 30322, USA

7 Cindy Xiong , University of Massachusetts - Amherst, Amherst, MA, 01003, USA

8 Yea-Seul Kim , University of Wisconsin - Madison, Madison, WI, 53706, USA

9 *In this Viewpoint article, we describe the persistent tensions between various
10 camps on the "right" way to conduct evaluations in visualization. Visualization as a
11 field is the amalgamation of cognitive and perceptual sciences and computer
12 graphics, among others. As a result, the relatively disjointed lineages in
13 visualization understandably approach the topic of evaluation very differently. It is
14 both a blessing and a curse to our field. It is a blessing, because the collaboration of
15 diverse perspectives is the breeding ground of innovation. Yet it is a curse, because
16 as a community, we have yet to resolve an appreciation for differing perspectives
17 on the topic of evaluation. We explicate these differing expectations and
18 conventions to appreciate the spectrum of evaluation design decisions. We
19 describe some guiding questions that researchers may consider when designing
20 evaluations to navigate differing readers' evaluation expectations.*

21 **I**magine that you are a visualization researcher
22 (skip this step if you already are one). You just got
23 the reviews back for your most recent submission.
24 Do you dare look?

25 Brave fictional researcher #1 Mojo opened the web
26 page containing the reviews. She has just submitted a
27 paper on how people perceive pie charts. The paper
28 contains three carefully designed studies that involve
29 showing participants pie charts with varying design
30 features to evaluate how quickly people can extract
31 key statistics from them. The paper was brutally
32 rejected. The reviews pointed out that the datasets
33 used to generate these pie charts were too limited,
34 and the study itself was too abstract and artificial.

35 Brave fictional researcher #2 Jojo also reads her
36 reviews. Her paper introduced a technique integrated
37 into a system, whose evaluation with in-lab partici-
38 pants assessed the technique's performance to

promote reflection of unconscious decision-making
39 strategies. Also brutally rejected, reviewers critiqued
40 the lack of control in the study and the abundance of
41 potential confounds. 42

43 Do these experiences sound familiar? As a research
44 community, we often struggle to design experiments
45 that balance readers' expectations. Maybe there is not
46 much we can do about the late hours we work, but there
47 must be something we can do about designing our stud-
48 ies in a way that ensures the quality of our work while at
49 the same time meeting these differing expectations.

50 We assert that this tension arises as a result of
51 *opposing lineages* in the visualization community. In
52 particular, one must look to the diversity of fields that
53 blend to form our field. VIS is influenced by fields such
54 as computer graphics, semiotics, HCI, cognitive science,
55 vision science, graphic design, and cartography, among
56 others. Each of these fields provides varying perspec-
57 tives and approaches to the challenges in visualization
58 research,¹ and, in particular, the preferred approach to
59 evaluation. For instance, researchers from a computa-
60 tion background tend to make contributions, such as
61 models or systems, which are often evaluated with dis-
62 crete metrics (e.g., accuracy) via simulation; HCI

63 researchers make contributions that are theoretical,
64 empirical, or artifacts, but evaluation typically involves
65 more human-centric experimentation using mixed quali-
66 tative and quantitative methods²; psychology research-
67 ers often make contributions in the form of empirical
68 findings from highly controlled laboratory experiments.
69 With this breadth and diversity in backgrounds and
70 expectations for evaluation, it is understandable that
71 there is often disagreement in visualization about what
72 constitutes a well-designed evaluation.

73 *AS A RESEARCH COMMUNITY, WE*
74 *OFTEN STRUGGLE TO DESIGN*
75 *EXPERIMENTS THAT BALANCE*
76 *READERS' EXPECTATIONS.*

77 We acknowledge that critique and disagreement
78 are a natural part of research. Particularly in the VIS
79 community, with such a broad range of backgrounds
80 and experiences that forged our field, the spectrum of
81 evaluation methods is likewise broad. In this View-
82 point article, we describe some of the lineages in our
83 community and their respective traditional expecta-
84 tions and norms for evaluation. We waded through the
85 collective amassment of our rejected evaluation
86 papers to bring you lovingly distilled lessons learned.
87 We focus on describing several decision points of eval-
88 uation design and suggest concrete guidance for ways
89 to think through these choices. This is not a guide for
90 how to pander to reviewers; rather, we hope that this
91 will help guide researchers through the often conflict-
92 ing goals of an evaluation.

93 BACKGROUND

94 Evaluation in VIS has been a hotly contested topic for
95 quite some time. Critiques often stem from the com-
96 munity's dissatisfaction with insufficient evaluation
97 techniques as well as lack of clarity or efficacy
98 in applying evaluation techniques. For instance,
99 researchers often express concerns with methods,
100 such as measuring the accuracy and time for users to
101 perform benchmark tasks with a visualization, which
102 does not provide researchers or developers insights
103 into the benefits of a visualization or visualization tool
104 nor actionable items to improve them. This reflects
105 a growing need in the visualization community to
106 consider evaluation techniques across all stages of
107 development to generate research-driven evidence
108 demonstrating the benefits of visualization.³

Emerging Work 109

110 In 2006, the BELIV workshop emerged as a respected
111 venue for novel evaluation methodologies.⁴ The initia-
112 tive has inspired an abundance of progress in evalua-
113 tion methodologies. For example, Shneiderman and
114 Plaisant proposed a method called multidimensional
115 in-depth long-term case studies (MILCS), which
116 assesses visualization tools based on observations,
117 interviews, surveys, and an expert user's likelihood to
118 achieve their goals with the tool over an extended
119 period of time to obtain multiple perspective on the
120 tool's effectiveness.⁵

121 Stasko⁶ proposed a framework describing the value
122 of visualization, which contained components describing
123 how a visualization can provide time savings and
124 insights, convey data, and inspire user confidence in
125 data. Inspired by this framework, Wall *et al.*⁷ created a
126 heuristic-based methodology that enables evaluators to
127 identify the strengths and weaknesses of a visualization
128 by quantitatively rating [heuristics](#), along several dimen-
129 sions. Researchers have also referenced evaluation tax-
130 onomies from education and humanities literature, such
131 as evaluating a visualization using learning outcomes
132 and Bloom's taxonomy.⁸

133 The initiative has also motivated researchers to cre-
134 ate guides on how to conduct evaluation studies. For
135 example, Elliott *et al.*⁹ introduced a lexicon of experimen-
136 tal design for empirical user studies, applying methodolo-
137 gies from human visual perception studies to evaluate
138 visualizations, describing novel experimental paradigms,
139 and dependent measures specific to the visualization
140 community. Sedlmair *et al.*¹⁰ provided a guide for con-
141 ducting design studies over nine stages (learn, winnow,
142 cast, discover, design, implement, deploy, reflect, and
143 write) and discussed potential pitfalls.

144 Other researchers categorized these evaluation
145 methods based on their high-level purposes in answering
146 a research question¹¹ and mapped out how these pur-
147 poses connect to the appropriate broader research con-
148 tribution.² For example, Munzner¹² proposed a nested
149 model that guides visualization researchers to select the
150 appropriate evaluation approach in four levels of visuali-
151 zation design and validation: characterizing the task and
152 data, abstracting the characterization into operations,
153 designing visual encoding and interactions, and creating
154 algorithms to execute the techniques efficiently.

155 **But** despite these efforts, the debate on how a
156 visualization should be evaluated rages on. Research-
157 ers and practitioners debate the criteria for measuring
158 the value of a visualization,⁶ expressing concerns on
159 the reproducibility of evaluation studies,¹³ and con-
160 tinue to share new expectations for visualization eval-
161 uations. Also, quietly in the background, lurks the

162 same debate in the form of a long discussion between
163 paper reviewers and paper chairs.

164 Contributions

165 Recently, the visualization research community's flag-
166 ship conference venue, IEEE VIS, has gone through a
167 remodel where, instead of separating submissions
168 into three subconferences (VAST, InfoVis, and SciVis),
169 submissions are now distributed into the following six
170 areas: theoretical and empirical; applications; systems
171 and rendering; representations and interaction; data
172 transformations; and analytics and decisions. These
173 areas may be roughly conceived of as contribution
174 types. Alternatively, Wobbrock and Kientz² describe
175 seven contribution types in HCI research, including
176 empirical (new observations), artifact (new tools),
177 methodological (new practices), theoretical (new con-
178 cept or model), dataset (new corpus), survey (new
179 reflection on a collection of past work), and opinion
180 (new perspective).

181 What makes visualization papers (or generally
182 speaking, any interdisciplinary academic research)
183 especially strong and unique contributions to the sci-
184 entific community is that one visualization paper often
185 touches on multiple areas and brings multiple forms of
186 contribution. For example, a visualization paper sub-
187 mitted to the theoretical and empirical track at IEEE
188 VIS could, in addition to its empirical contributions,
189 introduce a novel research method, a new dataset,
190 and a comprehensive background section that synthe-
191 sized a large amount of past work to be considered a
192 survey contribution. A paper submitted to the analyt-
193 ics and decisions area might contribute to the visuali-
194 zation community a new data analytic system (an
195 artifact), along with an empirical observation that
196 reveals new theoretical insights.

197 There are many ways to describe possible contribu-
198 tion types, but here we will focus on the following three:
199 factor, system, and technique, which often have different
200 evaluation expectations. A *factor* contribution usually
201 tells the story of how one design element impacts
202 the visualization and its interpretation. For example,
203 Ceja *et al.*¹⁴ demonstrated that the aspect ratio of a visu-
204 alization can influence how accurately people perceive
205 data. A *system* contribution showcases a novel tool to
206 help people build visualizations or analyze data, such as
207 Gratzl *et al.*'s work,¹⁵ which supports visual exploration of
208 rank data. For systems, evaluations are needed in order
209 to make claims about how effective the system is. Finally,
210 a *technique* contribution points to one specific compo-
211 nent (often in a system) and demonstrates that manipu-
212 lation of that technique can impact how people make

sense of visualizations. For example, Wall *et al.*¹⁶ demon- 213
strated an approach to displaying user interaction history 214
that may increase awareness of cognitive or societal 215
biases that drive behavior and decisions in data analytics. 216

This categorization complements the nested model 217
proposed by Munzner¹² by focusing on the end-product 218
of visualization research, the evaluation of which can 219
consist of any combination of the four levels from Munz- 220
ner's work.¹² 221

Challenges 222

223 While the multifaceted contributions of visualization
224 papers can lead to significant innovation, they also
225 introduce many challenges in the paper review pro-
226 cess. To evaluate a paper, the primary reviewer needs
227 to gather a group of reviewers with diverse back-
228 grounds and expertise to ensure a holistic evaluation
229 of the paper's contributions.

230 This is especially beneficial for papers that touch on
231 multiple areas and make multiple different contribu-
232 tions. However, one prominent issue often arises:
233 reviewers from different areas might judge the paper in
234 terms of its contribution in the one area they are familiar
235 with, without considering the other forms of contribu-
236 tion the paper brings. This can lead reviewers to find the
237 work underwhelming.

238 As a result, it becomes increasingly difficult for one
239 paper to reconcile the differing expectations from a group
240 of authors and a group of reviewers with diverse experien-
241 ces. The authors may feel pressured to make artificial
242 additions or omissions to please the reviewers. For exam-
243 ple, the long-running academic cliché laments that
244 reviewer number two makes unreasonable demands,
245 such as asking authors to conduct a full-fledged con-
246 trolled study, in stark contrast to other reviewers who
247 would prefer to see an ecologically valid study!

248 There are occasions where it is unnecessary or
249 impossible to run a perfectly balanced experiment or
250 user study that covers all possible confounds and
251 simultaneously maintains ecological validity. In fact,
252 many academics argue that it is a detriment to theo-
253 retical advancement to attempt to maximize external
254 validity in a given experiment.¹⁷ We need to collec-
255 tively acknowledge that no perfect experiment exists,
256 and one paper typically cannot solve an entire prob-
257 lem space. Papers ought to be judged on the experi-
258 ments conducted and contributions made, rather
259 than the ones they did not.

THINKING ABOUT STUDY DESIGN 260

261 We structure our discussion of evaluation guidance
262 around the three common types of contributions

RESEARCH DESIGN RESOURCES

There are a variety of articles on research methodologies in computing and psychology that we found helpful in building this guide to evaluation design. Listed below are a few of our favorites.

R. Elio, J. Hoover, I. Nikolaidis, M. Salavatipour, L. Stewart, and K. Wong, "About computing science research methodology," 2011.

L. Berkowitz and E. Donnerstein, "External validity is more than skin deep: Some answers to criticisms of laboratory experiments," *Amer. Psychologist*, vol. 37, no. 3, p. 245, 1982.

E. Wall, M. Agnihotri, L. Matzen, K. Divis, M. Haass, A. Endert, and J. Stasko, "A heuristic approach to value-driven evaluation of visualizations," *IEEE Trans. Vis. Comput. Graphics*, vol. 25, no. 1, pp. 491–500, 2018.

B. Shneiderman and C. Plaisant, "Strategies for evaluating information visualization tools: multi-dimensional in-depth long-term case studies," in *Proc. AVI Workshop BEyond Time Errors: Novel Eval. Methods Inf. Vis.*, 2006, pp. 1–7.

A. Burns, C. Xiong, S. Franconeri, A. Cairo, and N. Mahyar, "How to evaluate data visualizations across different levels of understanding," *IEEE Workshop Eval. Beyond-Methodol. Approaches Vis.*, pp. 19–28, 2020.

M. Sedlmair, M. Meyer, and T. Munzner, "Design study methodology: Reflections from the trenches and the stacks," *IEEE Trans. Vis. Comput. Graphics*, vol. 18, no. 12, pp. 2431–2440, 2012.

263 described earlier in the "Background" section: factor,
264 system, and technique contributions. How a user study
265 is designed by researchers and evaluated by reviewers
266 should depend on the type of contribution; the paper
267 claims to make. We identify a guide of several compo-
268 nents to help researchers design their studies and like-
269 wise help reviewers evaluate these studies.

270 Formulating Research Questions

271 Research begins by defining research questions and
272 corresponding claims researchers hope to address.
273 The research question should be directly connected
274 to the type of contributions; the paper sets out to
275 make. Munzner's¹² nested model for visualization
276 design and validation emphasized the importance of
277 asking the right research question, because address-
278 ing "the wrong problem" threatens the validity of every
279 step downstream in the research process. In this sec-
280 tion, we provide guidance on choosing an evaluation
281 that aligns with the specific research questions for a
282 feature, system, and technique.

283 *Example:* Let's imagine you are a visualization
284 researcher specializing in cognitive bias in visual analyt-
285 ics. You want to conduct a study that focuses on bias
286 mitigation. Inspired by previous work demonstrating
287 that having a user externalize their prior belief through
288 drawing can increase data recall,¹⁸ you come up with
289 the idea that people will be less susceptible to cognitive
290 biases in visual analytics if they can compare their

291 mental representation of the relationship between vari-
292 ables to the actual relationship between variables.
293 There are three ways you can approach your work: from
294 a factor, system, or technique perspective.

295 If you want to make a *factor* contribution and
296 determine possible associations or causal relations
297 between factors, it is probably a good idea to conduct
298 a highly controlled user study where the only differ-
299 ence between the conditions tested is that factor. For
300 the case study described, a good research question
301 might be "how does comparing a mental representa-
302 tion of the relationship between two variables to the
303 actual relationship influence one's interpretation
304 of data?"

305 If you want to make a *system* contribution, your
306 study should test whether your system improves an
307 existing visualization workflow, based on valid user
308 behaviors and intentions.^{6,12} The comparison to be
309 made here should be between the outcome from
310 when people use your system and the outcome from
311 when people do not. In the case study described, you
312 will probably want to design and build a visualization
313 system that can support bias mitigation. A good
314 research question might be "will people be less biased
315 when they analyze data using my system?" Notice
316 that the research question does not specifically talk
317 about the effect of any specific design elements in
318 your system, as it is up to you to how you want to
319 operationalize these elements.¹⁹ 319

TABLE 1. Description of three contribution types that we will focus on.

Contribution type	Description	Example RQ
Factor	Usually tells the story of how one design element impacts the visualization and its interpretation	How does comparing a mental representation of the relationship between two variables to the actual relationship influence one's interpretation of data?
System	Showcases a novel tool to help people build visualizations or analyze data	Will people be less biased when they analyze data using my system?
Technique	Points to one specific physical component (often in a system) and demonstrates that a manipulation of that technique can impact how people make sense of visualizations	Is the process of comparing mental representations to actual data mitigating cognitive bias in data analysis with my system?

320 Your system might include a novel technique that
 321 allows users to compare their mental representations
 322 of a data relationship, and you are sure that is the key
 323 to mitigate biases. It may be very tempting to add in
 324 your research question that this feature in the system
 325 mitigates biases. However, that turns your paper's
 326 contribution into a *technique*, rather than a *system*,
 327 and it will need to be evaluated differently, because a
 328 system is a complex collections of multiple techni-
 329 ques. Perhaps there is one technique in the system
 330 that is the key driver to mitigate bias, or perhaps the
 331 system pushes people to do analytic tasks in a certain
 332 order, and that order is what truly mitigates biases. If
 333 you additionally want to make claims about *why* your
 334 system works, you need to ask additional research
 335 questions at the technique level.

336 *DUE TO OPPOSING LINEAGES, THERE*
 337 *MAY BE A TENSION AMONG READERS'*
 338 *EXPECTATIONS FOR WHAT THE*
 339 *RESEARCH QUESTIONS OUGHT TO BE*
 340 *THAT YOU ADDRESS.*

341 If you want to make contributions at a *technique*
 342 level, you should think about how your technique can
 343 make a visualization system "better." The comparison
 344 you want to make in your study should be between
 345 the outcome for when people use a system with the
 346 target technique and the outcome for when people
 347 use the same system without the target technique,
 348 where other potential confounds are isolated. Notice
 349 how the system is kept constant in the comparison in
 350 this research question. This is because you need to
 351 justify the capability of the *technique itself* to combat
 352 the threat to the validity.¹² In the running example on

bias mitigation, you might want to make a claim of 353
 why your system works to mitigate bias. A good ques- 354
 tion might be "is the process of comparing mental rep- 355
 resentations to actual data mitigating cognitive bias 356
 in data analysis with my system?" These contribution 357
 types are tabulated in Table 1. 358

Due to opposing lineages, there may be a tension 359
 among readers' expectations for what the research 360
 questions ought to be that you address. It is possible 361
 then, or perhaps encouraged, to ask multiple research 362
 questions from different perspectives in your paper. A 363
 research question can be exploratory (which aims to 364
 navigate problem spaces to formulate hypotheses, often 365
 formulated prior to seeing data) or confirmatory (which 366
 aims to test a preexisting hypothesis the researchers 367
 have, often formulated after seeing data). It is critical to 368
 make sure, as mentioned in Munzner's work¹² that the 369
 research question always matches the output. The 370
 incongruency between the two tends to be a common 371
 source of critique from readers. It is on you to appropri- 372
 ately scope the research question and to motivate the 373
 problem space in the introduction and throughout this 374
 article to communicate why the exact questions contrib- 375
 ute to visualization design and systems. 376

Designing Conditions 377

Once you have formulated your research question, 378
 you should have a sense of what type of comparison 379
 you want to make in your study to validate your 380
 hypotheses and test the capabilities of your technique 381
 or system. This means coming up with the right test- 382
 ing conditions to answer your research questions.¹⁹ 383

Example: Let's continue with the bias mitigation 384
 example. Let's say your focus is at a *factor*-level and 385
 your research question is "how does comparing a 386
 mental representation of the relationship between 387
 two variables to the actual relationship influence 388
 one's interpretation of data?" The key comparison 389

390 here is how people interpret data in two scenarios:
 391 when they are able to compare a mental representa-
 392 tion of the relationship to the actual data, and when
 393 they are not able to (also known as the control condi-
 394 tion). Your study design should cover at least these
 395 two situations.

396 If your research focuses on evaluating your *system*,
 397 let's say your research question is "will people be less
 398 biased when they analyze data using my system?" You
 399 should include conditions that help you make the
 400 comparison between people's performance using your
 401 system versus another system or no system.

402 Let's say your contribution is at the *technique*-level,
 403 and your paper asks "is the process of comparing men-
 404 tal representations to actual data mitigating cognitive
 405 bias in data analysis with my system?" You should mini-
 406 mally test people's performance in your system with or
 407 without this technique by keeping everything else con-
 408 stant so you know exactly to what extent this tech-
 409 nique has an effect on user performance.

410 *THE MOST COMMON TENSION*
 411 *AMONG REVIEWERS WITH RESPECT*
 412 *TO EXPERIMENTAL CONDITIONS*
 413 *TYPICALLY LIES IN AN ISOLATION OF*
 414 *CONFOUNDING VARIABLES.*

415 The most common tension among reviewers with
 416 respect to experimental conditions typically lies in an
 417 isolation of confounding variables. To summarize, a
 418 system-level question where the conditions are "using
 419 system" and "not using system" can only warrant sys-
 420 tem-level conclusions, such as "we demonstrate that
 421 our system can help people perform better than no
 422 system." In this scenario, you cannot make technique-
 423 level claims and say "we demonstrate that our system
 424 can help people perform better because of this tech-
 425 nique" because you did not design study conditions to
 426 specifically test the effect of the technique, nor can
 427 you make factor claims about the mechanism of men-
 428 tal comparison toward mitigating bias since confound-
 429 ing variables were not isolated.

430 Internal and External Validity

431 In designing the experiment, you may wish to also
 432 ensure both internal and external validity.¹⁹ This forms
 433 another source of tension for authors to manage. One
 434 common complaint from reviewers for factor-based
 435 investigations is that these controlled studies lack

external validity, meaning how well the outcome of a 436
 study can be expected to apply to other settings in 437
 general. On the other hand, reviewers also often com- 438
 plain that system-based studies lack internal validity, 439
 which focuses on eliminating alternative explanations 440
 for a finding. So to address these common issues, we 441
 discuss a few study designs to provide an idea how to 442
 enhance both external and internal validity. Note, 443
 however, that it may not always be desirable to bal- 444
 ance both internal and external validity. For instance, 445
 theoretical advancements may be slowed by overin- 446
 dexing on external validity. 447

Example: For a *factor*-level investigation, say the 448
 question is "how does comparing a mental representa- 449
 tion of the relationship between two variables to the 450
 actual relationship influence one's interpretation of 451
 data." You probably want to operationalize all the rele- 452
 vant factors mentioned in your research question, and 453
 create sets of conditions to test the effect of each factor 454
 so you can pinpoint the factor(s) driving your effect. 455

For example, what is a mental representation of a 456
 relationship? Does this mean thinking about it? Drawing 457
 it out? Verbally describing *via* a sentence? What kind of 458
 variables do you want to focus on? Continuous varia- 459
 bles? Discrete variables? If you want to come up with 460
 generalizable results, you might want to test all types of 461
 mental representations on both continuous and discrete 462
 variables. Then, your conditions should cover all permu- 463
 tations of these two factors ({thinking, drawing, verbaliz- 464
 ing} x {discrete, continuous}) to yield six conditions. 465

But that is not all, you can keep asking yourself to 466
 further operationalize other factors in your research 467
 question: what is the actual relationship people should 468
 be comparing their mental representation with? Is this 469
 a visualization made from the underlying data? Per- 470
 haps a verbal description of a key insight? How would 471
 you measure "influence"? What about interpretation 472
 of data? Now you realize that this design space will 473
 expand exponentially as the number of permutations 474
 and combinations grows by adding more conditions to 475
 make your results more generalizable, but the amount 476
 of resources is finite. You cannot possibly run a well- 477
 powered study with 147 conditions and cram your find- 478
 ings in a nine-page paper. So now what? 479

We recommend you start by listing the entire 480
 design space of the experiment. Suppose you want to 481
 design an experiment with three experimental varia- 482
 bles (A, B, and C). First, you should list all the levels 483
 within each variable to exhaust the possibilities. For 484
 example, you could identify three levels for each vari- 485
 able (A1, A2, A3, B1, B2, B3, C1, C2, C3). Ideally, you can 486
 test all conjugated conditions, totaling 27 conditions. 487
 However, there are several reasons why testing all the 488

489 conditions may not be necessary. For example, prior
 490 work demonstrated that the combination of A1, B1,
 491 and C1 does not improve the task performance. You
 492 can eliminate the condition to save some time. You
 493 can also consult existing theories to narrow down the
 494 condition space.

495 *YOU CANNOT POSSIBLY RUN A WELL-*
 496 *POWERED STUDY WITH 147*
 497 *CONDITIONS AND CRAM YOUR*
 498 *FINDINGS IN A NINE-PAGE PAPER. SO*
 499 *NOW WHAT?*

500 Explicitly listing all the variables, levels, and the
 501 combined conditions is a necessary step toward think-
 502 ing about the entire space first to ensure internal
 503 validity. Researchers can then narrow down the space
 504 based on prior work and *communicate this process in*
 505 *their article*. This will enable reviewers and readers to
 506 follow the reasoning behind why a limited set of condi-
 507 tions have been tested and the rationale for those
 508 choices. Based on this, you can circle back to your
 509 research questions and re-scope it to match your
 510 study conditions (e.g., if you only tested the effect on
 511 continuous variables but not discrete variables, you
 512 may scope the research question down to explicitly
 513 focus on continuous variables). In a similar way, for
 514 system-level investigations, the goal is to demonstrate
 515 that your system actually works. To ensure internal
 516 validity, you want to make sure the only difference
 517 between your two conditions is whether your system
 518 is being used to complete the task or not. That means
 519 external variables like the task being tested, the par-
 520 ticipants' level of expertise in the task domain, among
 521 other things, should be kept constant.

522 If you also want to make *technique-level* claims,
 523 then you need to ensure that the only thing that
 524 differs between your conditions is whether that spe-
 525 cific technique exists or not. You should not compare
 526 a system with the targeted technique to a system
 527 without it, unless the two systems are identical to
 528 each other. Otherwise it violates internal validity; since
 529 there could be other differences in the system that
 530 make people perform better/worse, in addition to the
 531 technique of interest.

532 Choosing a Task

533 Now that you have your research questions formu-
 534 lated and your study conditions scoped, it is time to

think about what tasks you want users to complete 535
 for your study. The internal and external validity as 536
 well as the claims that you want to make should be 537
 considered in choosing a task in your study. 538

NOW THAT YOU HAVE YOUR 539
RESEARCH QUESTIONS FORMULATED 540
AND YOUR STUDY CONDITIONS 541
SCOPED, IT IS TIME TO THINK ABOUT 542
WHAT TASKS YOU WANT USERS TO 543
COMPLETE FOR YOUR STUDY. 544

For a *factor-level* contribution, it is important to 545
 abstract your task so that you can create an iso- 546
 lated environment to pinpoint the effect of your 547
 manipulation, such as how Amar *et al.*²⁰ and 548
 Brehmer and Munzner²¹ have abstracted analytic 549
 tasks for researchers to use to evaluate visualiza- 550
 tions. However, tension can arise when the task is 551
 too abstracted since the study can lose generaliz- 552
 ability to real-world settings. In these situations, the 553
 task may feel artificial to users, and the decisions 554
 they make in completing the task may no longer be 555
 good approximations of their actions in the real 556
 world. In these cases, an abundance of clarity in 557
 communicating the choices and tradeoffs can pre- 558
 empt many readers' concerns. 559

Example: In the bias-mitigation scenario, let's say 560
 we need some task that can capture a user's prior 561
 belief of a relationship, so we can see how that belief 562
 changes as biases are introduced or mitigated. Since 563
 you cannot randomly assign people to suddenly 564
 believe in one thing or another, if you want full control 565
 over people's prior beliefs, you likely need to make up 566
 an entirely artificial scenario and prime people with 567
 beliefs (e.g., the likelihood of a plant growing on an 568
 alien planet). But this will make the task seem con- 569
 trived, and participants might not take it seriously or 570
 respond in ways that would reflect their behaviors 571
 with real long-held beliefs. 572

Alternatively, you can take the organic approach 573
 and select scenarios where people's true belief can be 574
 easily measured and predicted and recruit people that 575
 hold a specific belief. This will likely resemble real- 576
 world scenarios more. However, beliefs that are easily 577
 measured and predicted are often associated with 578
 strong emotions, such as political orientation. For 579
 such strong or emotionally charged beliefs, you might 580
 not be able to observe changes in belief related to 581

582 biases. Although realistic, these emotional attach- 633
 583 ments may influence your results. 634

584 For a *system-* or *technique-*level contribution, 635
 585 abstracting a task to be short and simple can isolate 636
 586 confounding variables and make the specific task out- 637
 587 come more easily measured and controlled. However, 638
 588 systems are rarely designed for very simple tasks (e.g., 639
 589 Wall *et al.*'s work¹⁶). This makes evaluating performance 640
 590 via a simple task artificial and less useful for real-world 641
 591 usage scenarios. On the other hand, using real-world 642
 592 tasks usually means factors that may not be of research 643
 593 interest also play a part in the user workflow. This makes 644
 594 the outcome noisier to measure and the effect of a sys- 645
 595 tem or a technique enhancing performance on one spe- 646
 596 cific task more difficult to isolate from the influence of 647
 597 other factors or steps in the workflow. For example, let's 648
 598 say the system we designed mitigates bias by helping 649
 599 people see correlations in data more accurately. The 650
 600 abstracted control task may be to have users extract 651
 601 correlations from visualizations. In this case, the user 652
 602 views a visualization in the system and estimates a cor- 653
 603 relation value. These correlation estimates are com- 654
 604 pared to estimates made when the users view the same 655
 605 visualization in a different system. But reading correla- 656
 606 tion values from scatterplots is rarely the ultimate goal 657
 607 in a real-world system that supports a data analysis 658
 608 workflow. Just because people can more accurately 659
 609 read correlations from one system over another does 660
 610 not mean they are going to analyze the data, think about 661
 611 the data, or present the data in a less biased way. So 662
 612 alternatively, you may want to design a task that more 663
 613 closely resembles the real world.

614 Due to this tradeoff and potential tensions that 664
 615 can result between an abstracted, fully controlled task 665
 616 and a complex, real-world task, we recommend 666
 617 researchers to *consider planning for multiple studies*. 667
 618 It is possible to start with a more controlled setting in 668
 619 the first study to detect and quantify an effect; and 669
 620 then move to a more realistic setting in subsequent 670
 621 experiments where you examine whether the effect 671
 622 generalizes. This tradeoff in task choice is also a form 672
 623 of tradeoff of research contributions.

624 Picking a Dataset

625 The next consideration is with which dataset to design 673
 626 your visualization task. To ensure internal validity, you 674
 627 may want to generate your own datasets or look for 675
 628 datasets with specific characteristics. This way, you 676
 629 will have more control over what the visualization 677
 630 looks like and what type of analytic tasks users can 678
 631 perform. Some characteristics of a dataset that are 679
 632 manipulable might include the distribution (e.g.,

normal, uniform), how many discrete versus continu- 633
 ous variables there are, the number of abnormal/out- 634
 lier points, the size, etc. 635

Example: Continuing with the bias mitigation 636
 example, you want to see if people can become less 637
 biased in their data analysis with your intervention. 638
 Ideally you should test the effectiveness of your inter- 639
 vention with several different datasets with differing 640
 characteristics to see how much your results can gen- 641
 eralize before making claims. For example, it is possi- 642
 ble that your intervention can only mitigate biases 643
 when the dataset is normally distributed, ~~but~~ does not 644
 work when the dataset has more than 15% outliers. 645

646 However, there are too many possible characteris- 646
 647 tics of a dataset for it to be possible to control for 647
 everything. Identifying every manipulable characteris- 648
 tic of a dataset and creating a separate condition for 649
 each will likely consume too many resources, and 650
 dilute the focus of the research question. In addition, 651
 researcher-generated datasets may not resemble 652
 real-world datasets, reducing the external validity of 653
 the study as the study results may not generalize to 654
 real-world settings. These competing considerations 655
 can be another source of tension. A reader from a psy- 656
 chology background may expect these variables to be 657
 controlled for, while a reader from an HCI background 658
 may expect to see a realistic dataset. 659

660 We offer the following two considerations to help 660
 661 researchers who wish to balance these competing con- 661
 662 cerns in user studies: 1) start from a real-world dataset 662
 663 and manipulate the characteristics to gain control (e.g., 663
 664 by adding or removing columns or rows, altering the dis- 664
 665 tribution of a variable, etc.), or 2) list the assumptions of 665
 a realistic data generating process and simulate the 666
 data that also meets the control characteristics you 667
 need. These options give you the flexibility to have both 668
 realism in the user's perception of the data while main- 669
 taining control in the methods of analysis; but if you can- 670
 not have both, you need to choose one and stick to it. 671

672 From a Reviewer's Perspective

673 The flip-side of this guidance for researchers like- 673
 674 wise applies to reviewers. Reviewers should assess 674
 675 *the work according to how relevant the claims are* 675
 676 *to visualizations and whether the evidence supports* 676
 677 *the claims*. For papers that claim a factor-level con- 677
 678 tribution, assess how well the factor was isolated in 678
 679 influencing a phenomenon. For papers that claim a 679
 680 system-level contribution, assess whether the study 680
 681 is designed in such a way that it can capture differ- 681
 682 ences in alternative systems toward helping people 682
 683 achieve their goals under the same conditions. For a 683

684 paper that claims a technique-level contribution,
 685 reviewers should assess that the same system
 686 functions significantly differently with and without
 687 the technique. Furthermore, we encourage that
 688 reviewers assess *the contributions of the actual*
 689 *research that was done*, giving benefit of the doubt
 690 and, where appropriate, opportunity for authors to
 691 respond or revise when the language or framing of
 692 that research lacks precision (within reason).

693 CLOSING THOUGHTS

694 In an ideal world, we would be able to satisfy all
 695 internal and external validity goals in our evalua-
 696 tions. Realistically, however, we do not have infinite
 697 resources to realize all these oft-competing con-
 698 straints. In this Viewpoint, we have described some
 699 practical guidelines to help VisHikers navigate the
 700 galaxy of evaluation. While these guidelines will
 701 hopefully serve as a reasonable starting point in
 702 designing an evaluation and communicating those
 703 study design choices with precision, there are a
 704 number of other considerations we have barely
 705 touched on. We want to emphasize that there are
 706 other elements to consider in your study in addition
 707 to the experimental conditions.

708 *IMPORTANTLY, HOWEVER, THE*
 709 *PERFECT STUDY DOES NOT EXIST.*

710 For example, while it may be tempting to iden-
 711 tify multiple conditions to test in one study, unless
 712 you *conduct proper power analysis* to ensure you
 713 have the appropriate power to test your hypothe-
 714 ses, you risk collecting noisy measurements and
 715 observing unrepresentative effects. Similarly, you
 716 must *be careful not to overstate the contribution*.
 717 If you only compared people's performance using
 718 their system versus another state-of-the-art system
 719 using task A and B, then claim that "our system
 720 performs better at tasks A and B than the state-of-
 721 the-art system," rather than "our system performs
 722 better than this other system," or "our system is
 723 the best." Researchers must also consider how to
 724 carefully formulate their hypotheses, how to appro-
 725 priately measure the phenomena of interest, how
 726 to design a realistic data-generating process, and
 727 hopefully balance all of this within the context of a
 728 study that has some practical significance.

729 Importantly, however, the perfect study does not
 730 exist. There will always be tradeoffs that need to be

weighed and managed. We hope that this guidance 731
 will help re-enforce a critical and comprehensive lens 732
 for researchers to consider their evaluation designs. 733

ACKNOWLEDGMENTS 734

The authors would like to thank Jessica Hullman and 735
 John Stasko for their invaluable feedback. 736

REFERENCES 737

1. R. M. Kirby and M. Meyer, "Visualization collaborations: What works and why," *IEEE Comput. Graphics Appl.*, vol. 33, no. 6, pp. 82–88, Nov./Dec. 2013. 738
2. J. O. Wobbrock and J. A. Kientz, "Research contributions in human-computer interaction," *Interactions*, vol. 23, no. 3, pp. 38–44, 2016. 739
3. C. Plaisant, "The challenge of information visualization evaluation," in *Proc. Work. Conf. Adv. Vis. Interfaces*, 2004, pp. 109–116. 740
4. B. Lee, C. Plaisant, C. Parr, J. Fekete, and N. Henry, "Task taxonomy for graph visualization," in *Proc. AVI Workshop Beyond Time Errors: Novel Eval. Methods Inf. Vis.*, 2006, pp. 1–5. 741
5. B. Shneiderman and C. Plaisant, "Strategies for evaluating information visualization tools: Multi-dimensional in-depth long-term case studies," in *Proc. AVI Workshop Beyond Time Errors: Novel Eval. Methods Inf. Vis.*, 2006, pp. 1–7. 742
6. J. Stasko, "Value-driven evaluation of visualizations," in *Proc. 5th Workshop Beyond Time Errors: Novel Eval. Methods Vis.*, 2014, pp. 46–53. 743
7. E. Wall *et al.*, "A heuristic approach to value-driven evaluation of visualizations," *IEEE Trans. Vis. Comput. Graphics*, vol. 25, no. 1, pp. 491–500, Jan. 2019. 744
8. A. Burns, C. Xiong, S. Franconeri, A. Cairo, and N. Mahyar, "How to evaluate data visualizations across different levels of understanding," in *Proc. IEEE Workshop Eval. Beyond-Methodological Approaches Vis.*, 2020, pp. 19–28. 745
9. M. A. Elliott, C. Nothelfer, C. Xiong, and D. A. Szafir, "A design space of vision science methods for visualization research," *IEEE Trans. Vis. Comput. Graphics*, vol. 27, no. 2, pp. 1117–1127, Feb. 2021. 746
10. M. Sedlmair, M. Meyer, and T. Munzner, "Design study methodology: Reflections from the trenches and the stacks," *IEEE Trans. Vis. Comput. Graphics*, vol. 18, no. 12, pp. 2431–2440, Dec. 2012. 747
11. N. Elmqvist and J. S. Yi, "Patterns for visualization evaluation," *Inf. Vis.*, vol. 14, no. 3, pp. 250–269, 2015. 748
12. T. Munzner, "A nested model for visualization design and validation," *IEEE Trans. Vis. Comput. Graphics*, vol. 15, no. 6, pp. 921–928, Nov./Dec. 2009. 749

780 13. J.-D. Fekete and J. Freire, "Exploring reproducibility in
781 visualization," *IEEE Comput. Graphics Appl.*, vol. 40,
782 no. 5, pp. 108–119, Sep./Oct. 2020.

783 14. C. R. Ceja, C. M. McColeman, C. Xiong, and S. L. Fran-
784 coneri, "Truth or square: Aspect ratio biases recall of
785 position encodings," *IEEE Trans. Vis. Comput. Graphics*,
786 vol. 27, no. 2, pp. 1054–1062, Feb. 2021.

787 15. S. Gratzl, A. Lex, N. Gehlenborg, H. Pfister, and M. Streit,
788 "LineUp: Visual analysis of multi-attribute rankings,"
789 *IEEE Trans. Vis. Comput. Graphics*, vol. 19, no. 12,
790 pp. 2277–2286, Dec. 2013.

791 16. E. Wall, A. Narechania, A. Coscia, J. Paden, and A.
792 Endert, "Left, right, and gender: Exploring interaction
793 traces to mitigate human biases," *IEEE Trans. Vis.*
794 *Comput. Graphics*, vol. 28, no. 1, pp. 966–975, Jan. 2022.

795 17. B. J. Calder, L. W. Phillips, and A. M. Tybout, "The
796 concept of external validity," *J. Consum. Res.*, vol. 9,
797 no. 3, pp. 240–244, 1982.

798 18. Y.-S. Kim, K. Reinecke, and J. Hullman, "Explaining the
799 gap: Visualizing one's predictions improves recall and
800 comprehension of data," in *Proc. CHI Conf. Hum.*
801 *Factors Comput. Syst.*, 2017, pp. 1375–1386.

802 19. P. C. Cozby and S. Bates, *Methods in Behavioral*
803 *Research*. Mountain View, CA, USA: Mayfield Pub., 1985.

804 20. R. Amar, J. Eagan, and J. Stasko, "Low-level
805 components of analytic activity in information
806 visualization," in *Proc. IEEE Symp. Inf. Vis.*, 2005,
807 pp. 111–117.

808 21. M. Brehmer and T. Munzner, "A multi-level typology of
809 abstract visualization tasks," *IEEE Trans. Vis. Comput.*
810 *Graphics*, vol. 19, no. 12, pp. 2376–2385, Dec. 2013.

EMILY WALL is an Assistant Professor with Emory University, 811
Atlanta, GA, USA. Her research focuses on decision-making 812
with data and visualizations, particularly as applied to prob- 813
lems of societal concern. She received the Ph.D. degree in 814
computer science from the Georgia Institute of Technology, 815
Atlanta. She is the corresponding author of this article. Con- 816
tact her at emily.wall@emory.edu. 817

CINDY XIONG is an Assistant Professor with the University of 818
Massachusetts Amherst, Amherst, MA, USA. Her research 819
interests include perception, cognition, and data visualization, 820
and she investigates how humans perceive, interpret, and 821
make decisions from visualized data. She received the Ph.D. 822
degree in psychology from Northwestern University, Evanston, 823
IL, USA. Contact her at cindy.xiong@cs.umass.edu. 824

YEASEUL KIM is currently an Assistant Professor with the 825
University of Wisconsin-Madison, Madison, WI, USA. Her 826
research focuses on developing tools and algorithms to help 827
people with varying abilities interact with data and visualiza- 828
tions. She received the Ph.D. degree in information science 829
from the University of Washington, Seattle, WA, USA. Contact 830
her at yeaseul.kim@cs.wisc.edu. 831

Contact department editor Theresa-Marie Rhyne at 832
theresamarierhyne@gmail.com 833