

1 Machine learning-based evidence and attribution mapping of 100,000
2 climate impact studies

3
4 Max Callaghan^{1,2}, Carl-Friedrich Schleussner^{3,5}, Shruti Nath^{3,4}, Quentin Lejeune³, Thomas R. Knutson⁶,
5 Markus Reichstein^{7,8}, Gerrit Hansen⁹, Emily Theokritoff^{3,5}, Marina Andrijevic^{3,5}, Robert J. Brecha^{3,10},
6 Michael Hegarty³, Chelsea Jones³, Kaylin Lee³, Agathe Lucas, Nicole van Maanen^{3,5}, Inga Menke³, Peter
7 Pfliederer^{3,5}, Burcu Yesil³, Jan C. Minx^{1,2}

8
9 ¹ Mercator Research Institute on Global Commons and Climate Change, Berlin, Germany.

10 ² Priestley International Centre for Climate, University of Leeds, Leeds, LS2 9JT, UK.

11 ³ Climate Analytics, Berlin, Germany

12 ⁴ Institute of Atmospheric and Climate Sciences, ETH Zürich, Switzerland

13 ⁵ Integrative Research Institute on Transformations of Human-Environment Systems, Humboldt
14 University, Berlin, Germany

15 ⁶ NOAA/Geophysical Fluid Dynamics Laboratory, Princeton, NJ, 08540, USA

16 ⁷ Max Planck Institute for Biogeochemistry, Department Biogeochemical Integration, D-07701 Jena,
17 Germany

18 ⁸ Michael Stifel Center Jena for Data-driven and Simulation Science, Jena, Germany

19 ⁹ Robert Bosch Stiftung GmbH, Berlin, Germany

20 ¹⁰ Hanley Sustainability Institute, Renewable and Clean Energy Program and Physics Dept., University
21 of Dayton, Dayton, Ohio, USA

22

23 Accepted for publication in *Nature Climate Change*, version as of Aug. 31, 2021

24

25 **Abstract**

26 **Increasing evidence suggests that climate change impacts are already observed around the world.**
27 **Global environmental assessments face challenges to appraise the growing literature. Here we use**
28 **the language model BERT to identify and classify studies on observed climate impacts, producing a**
29 **comprehensive machine-learning-assisted evidence map. We estimate that 100,724 (62,950-162,838)**
30 **publications document a broad range of observed impacts. By combining our spatially resolved**
31 **database with grid cell level human-attributable changes in temperature and precipitation, we infer**
32 **that attributable impacts may be occurring across 80% of the world's land area where 85% of the**
33 **population reside. Our results reveal a substantial 'attribution gap' as robust levels of evidence for**
34 **attributable impacts is twice as prevalent in high income than low income countries. While gaps**
35 **remain on confidently establishing attributable climate impacts at the regional and sectoral level,**
36 **this database illustrates the potential current impact of anthropogenic climate change across the**
37 **globe.**

38

39

40 There is overwhelming evidence that the impacts of climate change are already being observed in human
41 and natural systems¹. These effects are emerging in a range of different systems and at different scales,
42 covering a broad range of research fields from glaciology to agricultural science, and marine biology to
43 migration and conflict research². The evidence base for observed climate impacts is expanding³, and the
44 wider climate literature is growing exponentially^{4,5}. Systematic reviews and systematic maps offer
45 structured ways to collectively identify and describe this evidence while maintaining transparency,
46 attempting to ensure comprehensiveness and reduce bias⁶. However, their scope is often confined to very
47 specific questions covering no more than dozens to hundreds of studies.

48
49 In the climate science community, evidence-based assessments of observed climate change impacts are
50 performed by the Intergovernmental Panel on Climate Change (IPCC)². Since the first Assessment Report
51 (AR) of the IPCC in 1990, we estimate that the number of studies relevant to observed climate impacts
52 published per year has increased by more than two orders of magnitude (Fig. 1a). Since the third AR,
53 published in 2001, the number has increased ten-fold. This exponential growth in peer-reviewed scientific
54 publications on climate change^{4,5} is already pushing manual expert assessments to their limits. To address
55 this issue, recent work has investigated ways to handle big literature in sustainability science by scaling
56 systematic review and map methods to large bodies of published research using technological innovations
57 and machine learning methods⁷⁻¹¹. Much of this work builds on a related literature that has applied natural
58 language processing techniques to problems of evidence synthesis in the health sciences¹²⁻¹⁴.

59
60 Fully utilising the available knowledge on emerging climate change impacts is key to informing global
61 policy processes¹⁵ as well as regional and local risk assessments and on-the-ground action on climate
62 adaptation^{16,17}. While the global policy process may be served well with literature assessments presenting
63 results aggregated on the level of continents or world regions^{2,18}, informing climate adaptation typically
64 requires more highly localised and contextualised information on climate impacts^{19,20}.

65
66 Another core challenge of literature reviews and assessments of observed climate impacts relates to the
67 question of whether climate impacts can be attributed to anthropogenic forcing²¹. While anthropogenic
68 climate change signals have been identified in observed trends in a number of variables²¹ including
69 temperature²², precipitation²³, sea level rise²⁴, or water resources²⁵, and selected extreme weather²⁶ events,
70 the confidence in these assessments is still subject to substantial regional variations and remains relatively
71 tentative at smaller spatial scales even if very high confidence levels can be reached for larger scale (e.g.,
72 global scale) attribution findings. Confidence also strongly depends on the variable being considered, and
73 specifically decreases further down the impact chain, i.e. for indicators of changes in human and natural
74 systems that are driven by changes in other climate impact variables²¹. In addition, methodological
75 approaches and robustness criteria for climate change attribution differ widely between studies and
76 disciplines, requiring expert judgement on a case-by-case basis in order to compile a comprehensive
77 evidence base.

78
79 This points towards the added value of joining the body of evidence documenting regional or local-scale
80 studies about climate impacts linked to common climate drivers such as temperature and precipitation
81 change to a spatially resolved detection/attribution database of those variables.

82

83 Using BERT, a state of the art deep learning language representation model²⁷, we develop a machine
84 learning pipeline to identify, locate and classify studies on observed climate impacts at a scale beyond
85 that which is possible manually (see Extended Figure 1). We combine this spatially resolved dataset with
86 an approach to attributing observed trends in surface temperature and precipitation at the grid cell level
87 (5° x 5° and 2.5° x 2.5° cells respectively) to human influence on the climate. In doing so, we establish a
88 new paradigm for assessing the impacts of climate change across human and natural systems.

91 **Mapping over 100,000 impact studies**

92
93 We searched two large bibliographic databases (Web of Science and Scopus) using an inclusive and
94 transparent search method to systematically identify the literature on climate impacts. We assessed
95 comprehensiveness by ensuring that our search string returned all references from tables 18.5-18.9 in
96 AR5 WGII, which deal with the detection and attribution of climate impacts. Recent breakthroughs in
97 natural language processing (NLP) have extended the capabilities of text classification. BERT
98 (Bidirectional Encoder Representations from Transformers) is a deep learning language model trained
99 using semi-supervised learning on massive corpora to represent text where word representations are
100 dependent on context. Such models are able to some extent able to capture the context-dependent
101 meanings of texts. The pretrained model can be fine-tuned on downstream tasks, and has achieved state of
102 the art results across a range of NLP tasks. Using training data assembled by collaboratively screening
103 and coding 2,629 abstracts, we use supervised machine learning, fine-tuning the smaller and faster BERT
104 variant DistilBERT²⁸, to classify, also based on the abstract text, documents relevant to understanding the
105 observed impacts of climate change in general, and to predict the human or natural systems for which
106 they document impacts (i.e., the impact categories), as well as the climate variable(s) driving the
107 documented impacts. Uncertainty estimates for the predictions are derived from bootstrapping. We
108 employ a nested cross-validation approach to hyperparameter tuning, model selection and classifier
109 evaluation, and find that our binary inclusion classifier achieves an average F1 score of 0.71, and ROC
110 AUC score of 0.92. The prediction of impact type is achieved with an average macro F1 score of 0.84
111 while the prediction of climate driver is achieved with an average F1 score of 0.79 (see Methods section
112 and Extended Figures 1-5 for a detailed explanation of the labelling, machine learning approach and
113 classifier performance).

114
115 Our query returned 601,677 unique documents (Fig. 1a): many more than would have been possible to
116 screen by hand. Of these we estimate that 102,160 (64,386-164,274) documents are relevant to
117 understanding the observed impacts of climate change in general, based on the spread of
118 inclusion/exclusion predictions obtained from our model via bootstrapping (Fig. 1a.). This base of
119 relevant publications has grown substantially through the IPCC assessment cycles. 46,442 (34,473-
120 87,861) articles have been published in the sixth assessment cycle so far; this represents more than twice
121 the number of studies published during the AR5 period.

122
123 We used a geoparser pre-trained using neural networks²⁹ to extract structured geographic information
124 from the titles and abstracts of the studies in our database. Although the number of relevant studies in
125 North America, Asia, and Europe is much higher than in South America, Africa, and Oceania, there is a
126 large body of relevant studies available on all continents (fig 1.c). Adjusted for population

127 (Supplementary Fig. 1), the number of papers focusing on Oceania far exceeds the size of the literature
128 devoted to other continents, with Africa and Asia receiving the least attention per million inhabitants. The
129 relevant publications are also unevenly distributed across impact categories, with by far the largest
130 number of studies 34,988 (18,520 - 65,666) documenting impacts on terrestrial and freshwater ecosystems
131 (Fig 1.b.). However, the category with the comparably smallest coverage--mountains, snow and ice--still
132 has 6,307 (3,526 - 12,228) studies.

133
134 In contrast to the map of observed impacts produced by the IPCC, we do not only include papers which
135 formally attribute impacts to observed trends in climate. Instead, we take a more comprehensive approach
136 reflecting that our objective is to map all possibly relevant studies on climate-related changes, rather than
137 a list of studies where the relationship between an observed climate trend and specific impacts has been
138 demonstrated with high confidence, or even linked to human influence on the climate. This includes
139 studies attributing impacts to observed trends in climate variables, even where the authors do not attribute
140 these trends to human influence, such as, for example, a study documenting the influence of the date of
141 snowmelt on the phenology and population growth of mammals³⁰. In addition, we include studies which
142 provide evidence on the sensitivity of human or natural systems to climate metrics, such as how heart
143 disease mortality responds to variations in temperature³¹. Finally, we include documents describing the
144 impacts of extreme events and studies which detect significant trends in climate variables or climate
145 extremes³², regardless of whether or not these trends are in line with the expected effects of
146 anthropogenic climate change. We exclude all studies which only describe potential or modelled impacts
147 of future climate change.

148 149 **Combining geolocated literature with climate information**

150 To add context on the role of anthropogenic climate change in driving impacts, or more precisely the role
151 of historical changes in anthropogenic climate forcing agents such as greenhouse gases and aerosols, we
152 combine our literature database of studies selected using machine learning with spatially explicit analysis
153 of detectable and attributable trends in two key climate variables. Combining evidence from climate
154 model simulations and observational datasets allows identification of trends likely attributable in part to
155 anthropogenic climate change for near-surface temperature and precipitation at the level of 5 degree
156 (temperature) or 2.5 degree (precipitation) grid cells^{22,23}. Here we apply this methodology to analyse
157 trends from 1951 to updated observational data until 2018 for temperature (Fig.2a) and until 2016 for
158 precipitation (Fig.2b). Grid cells in categories +2 or +3 show where trends cannot be explained by
159 internal variability and are either consistent with or greater than the expected change in climate model
160 simulations that include anthropogenic forcing agents. We infer that these cells display detectable and at
161 least partly attributable trends (see Methods for more details).

162
163 We next resolve the structured geographic information extracted from our studies, which range from
164 continental scale down to individual watersheds or communities, to sets of grid cells (Extended Fig. 9,
165 Methods). We can then derive the weighted number of studies per grid cell according to the number of
166 grid cells to which each study relates. By combining studies related to temperature or precipitation with
167 the gridded information on attributable trends in temperature and precipitation, this provides a necessary
168 (though not necessarily sufficient) condition for a systematic two-step attribution to anthropogenic
169 activities of the impacts predicted by the classifier³³. Where studies documenting impacts associated with
170 changes in temperature or precipitation co-occur with attributable trends in those variables, we claim that

171 there is at least preliminary evidence for attributable impacts in these areas. This approach is similar in
172 nature to the “joint attribution” applied in IPCC AR4^{34,35}.

173
174 In general, we note that this type of automated assessment procedure is no substitute for careful
175 assessment by experts, but can identify large numbers of studies for a region that may point toward
176 attributable human influence on impacts. Confidence in multi-step attribution claims depends on
177 confidence in the attribution of the individual components (steps) along with the confidence or limitation
178 in linking the different steps in the proposed causal chain³⁵. One limitation of the partially automated
179 two-step attribution approach is that we cannot verify that every temperature or precipitation trend cited
180 in impact studies matches, either in sign, magnitude or time period, those attributed to human influence
181 by the regional detection and attribution studies for temperature²² and precipitation²³. This is a greater
182 problem for studies driven by precipitation, where both wetting and drying trends occur with greater
183 temporal variation, though these make up the minority of partially attributed studies and grid cells. We
184 also note that not all studies in our database document impacts in response to trends in climate variables.
185 Where impacts are attributed to extreme events or variation in temperature or precipitation, the fact that
186 recent trends in temperature or precipitation can be attributed to human influence provides important
187 context, but does not allow robust attribution of those impacts. These factors limit confidence in our cases
188 of potential attribution of impacts to anthropogenic forcing. Our approach could be extended with more
189 fine-grained analysis of studies or with attribution of additional signals in climate variables in order to
190 make more robust attribution statements.

191
192 For 80% of global land area (excluding Antarctica), trends in temperature and/or precipitation can be
193 attributed at least in part to human influence on the climate (purple cells, Fig. 2c). Using gridded
194 population density data³⁶, we calculate that this covers 85% of the world’s population. The majority of
195 land grid cells show attributable warming trends, with exceptions where trends cannot be robustly
196 distinguished from internal variability (white cells, category 0) or where there is insufficient data to
197 establish trends (grey cells). For precipitation, attributable wetting and drying trends are found with
198 greater geographical variation. There are also more grid cells where a trend in precipitation cannot be
199 established, or where the observed trend is opposite in sign to that simulated by climate model historical
200 simulations (purple and yellow cells, +-4).

201
202 Though most of the world’s population resides in areas where trends in temperature and or precipitation
203 can be at least partially attributed to human influence, there is substantial geographical variation in the
204 degree to which the impacts of temperature and precipitation on human and natural systems have been
205 studied. We characterise areas with fewer than 5 weighted studies per grid cell as displaying low
206 evidence, areas with between 5 and 20 weighted studies as robust evidence, and areas with more than 20
207 weighted studies as high evidence.

208
209 For 48% of global land area (hosting 74% of global population), we find robust or high evidence of
210 impacts on human and natural systems colocated with attributable temperature or precipitation trends
211 (Fig. 2c). Areas with this combination of evidence are indicated by the darker purple cells. These
212 constitute almost all grid cells in Western Europe, North America, South and East Asia, and there are
213 parts of all continents which have similar pockets of substantial preliminary evidence.

214

215 However, for 33% of global land area (hosting 11% of global population), although there is evidence that
216 long-term trends in precipitation and temperature are attributable at least in part to human influence, there
217 is relatively little evidence in the existing literature about how these trends impact human and natural
218 systems (Fig. 2c lightest purple shading). This imbalance suggests, in line with research measuring
219 climate impacts using remote sensing³⁷, that the lack of evidence in individual studies is because these
220 locations are less intensively studied, rather than an absence of impacts in these areas. Parts of Western
221 Africa, South-east, Western and Northern Asia contain several light red grid cells where there is evidence
222 to suggest that the climate (temperature and/or precipitation) has changed because of human influence,
223 but there is little evidence on how this may be impacting human and natural systems. These demonstrable
224 evidence gaps suggest a lack of impacts research commensurate with current knowledge of how the local
225 climate (temperature and/or precipitation) is changing.

226
227 Some of the spatial features can be explained by the geographical characteristics. Among the regions with
228 limited evidence are vast, sparsely populated and difficult to reach areas with a comparable uniform
229 biosphere and climate such as Siberia or the Saharan desert. But beyond these features, our results clearly
230 reveal a substantial 'attribution gap'. We find that 23% of the population of low income countries live in
231 areas with low impact evidence despite at least partially attributable trends in temperature and/or
232 precipitation (Fig. 2.d). In high income countries, this figure is only 3%. A density of 5 studies per grid
233 cell or more with attributable impacts is 1.76 times as prevalent by population for high income countries
234 (88%) as for low income countries (50%), while a density of 20 studies or more with attributable impacts
235 is more than 4 times as prevalent (81% compared to 17%).

236
237 In the remaining grey grid cells (Fig. 2c), trends in precipitation and temperature have not been attributed
238 to human influence on the climate according to the methodology in refs. 18 and 19, as applied to CMIP6
239 models. This does not rule out the possibility that some trends in precipitation or temperature have
240 occurred in these regions that have been driven, at least in part, by human influence on the climate.
241 However, due to various factors, such as lack of adequate observational data, high levels of natural
242 variability compared to the climate change signal, or limitations in modelling or estimated climate
243 forcings, some observed changes that actually include anthropogenic contributions may not yet be
244 attributable at the grid cell level. This categorisation of individual gridpoints may well change as new
245 observational data are collected, as models improve, as the global climate continues to warm, or as
246 detection/attribution methodologies improve. Darker grey grid cells (10% of analyzed land area) indicate
247 where there are no detectable trends in temperature or precipitation that can be attributed to human
248 influence at a grid cell level, but where there nevertheless appears to be substantial evidence that local
249 trends in some climate variables lead to impacts on human and natural systems. For example, many
250 studies refer to the impacts of temperature in the state of Western Australia, but of the 40 grid cells in the
251 state, an attributable temperature trend can be demonstrated for 22 cells. For 16 of the remaining cells a
252 lack of data means that a detectable trend cannot be established, and for the remaining 2 cells, no
253 attributable trend can be established.

254
255 The lightest grey cells (17% of land area) describe areas where we do not detect anthropogenic influence
256 on regional temperature or precipitation and find few publications about the impacts of temperature or
257 precipitation on human and natural systems in those areas. Apart from high latitudes and over the ocean,
258 these cells are primarily in Africa. For example, in the light grey patch over the central part of sub-

259 Saharan Africa, either limitations of observed data, models, or low signal to noise imply that we are
260 unable to attribute temperature or precipitation trends to human influence on the climate using the
261 methodologies employed here (see extended fig. 4); further, we have identified few studies analysing the
262 impacts of climate change on human and natural systems in those regions. These evidence gaps constitute
263 significant blind spots in understanding of climate impacts, and in some cases understanding of
264 attributable anthropogenic influence on regional precipitation and/or temperature.

265
266 In total, 57,366 studies discuss impacts related to a driver which our analysis suggests can be attributed in
267 part to human influence on the climate in at least one grid cell to which the study refers. We find
268 hundreds of partially or mostly attributable studies (where there are attributable trends in the relevant
269 climate variable for at least 1% or more than 50% of grid cells respectively) in each impact category
270 across all continents (Fig. 3, indicated by the darker green and purple bars). This figure ranges from 268
271 (143-514) studies of impacts on mountains, snow and ice in Africa to 7,835 (4,308-13,552) studies of
272 impacts on terrestrial ecosystems in North America. Wide confidence intervals here reflect the compound
273 uncertainty deriving from classification of relevance, impact and driver.

274
275 Our analysis also allows quantification of how the share of research on each impact category varies from
276 continent to continent. For example, research on human and managed systems makes up 12% of all
277 research globally, but only 10% of research in Europe, compared to 19% in Africa. This focus on human
278 and managed systems in Africa is remarkable given that the absolute numbers of studies in Africa (1,466)
279 is similar to that in Europe (1,799) despite the vast difference in total numbers of studies between the two
280 continents. This greater share of research in Africa documents impacts in human and managed systems
281 may reflect the high vulnerability of particularly sub-Saharan Africa to climate impacts³⁸.

282
283
284

285 **Discussion and conclusion**

286 We develop a two-step attribution process which combines a transparent and reproducible^{39,40} machine
287 learning approach to identifying studies on observed climate impacts with model-based assessments of
288 detectable anthropogenic contributions to historical temperature and precipitation trends. Using machine
289 learning to scale up evidence synthesis allows us to map 100,000 studies of climate impacts, providing a
290 comprehensive picture of the evidence base. Bringing together these two lines of evidence on climate
291 change and climate impacts provides a new bridge between the climate science community and the
292 impacts, adaptation, and vulnerabilities communities, and highlights the synergistic nature of their
293 approaches.

294
295 Our spatially resolved approach allows for a systematic provision of regional to local, sector-specific
296 climate impact information to local or regional experts and adaptation practitioners. This offers
297 perspectives for a novel climate service supporting the uptake of scientific information in local contexts
298 and providing relevant information for adaptation action. Second, the quantification of an “attribution
299 gap” highlights the need for more research on climate impacts in low income countries. Furthermore, the
300 automated nature of the assessment allows for continuous updating of the database, creating a ‘living’

301 evidence map that can also be improved and extended by incorporating additional sources of relevant
302 publications (e.g. non-English speaking evidence, or improved/expanded regional detection/attribution
303 studies) and targeted assisted learning in regional or topical areas of interest.

304
305 The compiled database is vast, but neither complete nor perfect. Our systematic query-based literature
306 search is extensive, but will also exclude some relevant studies. The selection and categorisation of
307 studies was achieved using machine learning, meaning that results are subject to additional uncertainties,
308 which compound for each level of classification. Further, documents were coded only at the abstract
309 level, and only the abstracts were used as inputs to our classifiers. Given the relative simplicity of the type
310 of information we extract (focusing on the impact area studied and the documented driver), we expect
311 them to be covered in the abstract, which provides the condensed summary of the study's findings.
312 Applying classifiers to noisy full texts which contain contextual information and related research as well
313 as the results and topic of a study would greatly increase the risks of false positives. We thus find our
314 approach well justified for such high-level syntheses.

315
316 The database we assemble will also incorrectly exclude some relevant documents and contain some
317 documents that have been incorrectly included or incorrectly coded, but the approach enables us to report
318 both classifier performance and associated uncertainties. Additionally, some included studies may be of
319 low quality, as no process for critical appraisal (a key component of formal systematic reviews) was
320 followed either by human reviewers or in the machine learning pipeline. In the case of systems subject to
321 other anthropogenic interference such as the global biosphere, managed systems such as agriculture, or
322 human systems themselves, identifying a robust climate change driver requires careful assessment of
323 other socio-economic factors^{41,42}, adding additional levels of complexity⁴³.

324
325 The two-step attribution process is also only applied for the subset of papers which provide evidence on
326 impacts driven by temperature and precipitation. Exploring the role of human influence for studies
327 analysing the effects of factors other than trends in mean temperature or precipitation as the main driver
328 would require additional attribution strategies, but these could in principle be combined with individual
329 studies in similar ways. There is a growing literature on attributable human influence on a number of
330 climate metrics at the regional scale as well as extreme events⁴⁴⁻⁴⁶, and therefore much scope for
331 expansion of this approach. Finally, we note that plausible causal chains of cascading impacts are not
332 covered by our attribution approach (such as temperature driving an increase in drought, leading to
333 reduced agricultural yields) except where studies address each part of the causal chain.

334
335 These caveats highlight that the type of machine learning-assisted evidence map we present here is no
336 substitute for careful assessment by experts, either in the context of a gold-standard systematic review⁴⁷ or
337 in IPCC assessments. However, in an age of "big literature"^{7,9}, it is an invaluable complement. The use of
338 machine learning means we consider more evidence than would otherwise be feasible, showing where
339 evidence appears to be more prevalent and where important gaps can be observed. While traditional
340 assessments can offer relatively precise but incomplete pictures of the evidence, our machine-learning-
341 assisted approach generates an expansive preliminary but quantifiably uncertain map. Further, it enables
342 us to provide an automated, living systematic map of climate impacts that can be readily updated.
343 Ultimately, we hope that our global, living, automated, and multi-scale database will help to jump-start a
344 host of reviews of climate impacts on particular topics or particular geographic regions.

345
346 Machine-learning pipelines as developed here will be useful to prepare the IPCC for the age of big
347 literature by scaling systematic evidence mapping approaches. However, our results also show how
348 synthesis and transparency can be lifted to new levels by combining so-far disparate lines of evidence and
349 reporting classifier performance as well as associated uncertainties. If science advances by standing on
350 the shoulders of giants, in times of ever-expanding scientific literature giants' shoulders become harder to
351 reach. Our computer-assisted evidence mapping approach can offer a leg-up.
352

353 References

- 354
- 355 1. Cramer, W. *et al.* Detection and attribution of observed impacts. in *Climate Change 2014: Impacts,*
356 *Adaptation, and Vulnerability. Part A: Global and Sectoral Aspects. Contribution of Working Group*
357 *II to the Fifth Assessment Report of the Intergovernmental Panel of Climate Change* (eds. Field, C. B.
358 et al.) 979–1037 (Cambridge University Press, 2014).
 - 359 2. IPCC. *Climate Change 2014: Impacts, Adaptation, and Vulnerability. Part A: Global and Sectoral*
360 *Aspects. Contribution of Working Group II to the Fifth Assessment Report of the Intergovernmental*
361 *Panel on Climate Change.* (Cambridge University Press, 2014).
 - 362 3. Hansen, G. The evolution of the evidence base for observed impacts of climate change. *Curr. Opin.*
363 *Environ. Sustain.* **14**, 187–197 (2015).
 - 364 4. Haunschild, R., Bornmann, L. & Marx, W. Climate Change Research in View of Bibliometrics.
365 *PLOS ONE* **11**, e0160393 (2016).
 - 366 5. Grieneisen, M. L. & Zhang, M. The current status of climate change research. *Nat. Clim. Change* **1**,
367 72–73 (2011).
 - 368 6. Haddaway, N. R. & Pullin, A. S. The Policy Role of Systematic Reviews: Past, Present and Future.
369 *Springer Sci. Rev.* **2**, 179–183 (2014).
 - 370 7. Callaghan, M. W., Minx, J. C. & Forster, P. M. A topography of climate change research. *Nat. Clim.*
371 *Change* **10**, 118–123 (2020).

- 372 8. Porciello, J., Ivanina, M., Islam, M., Einarson, S. & Hirsh, H. Accelerating evidence-informed
373 decision-making for the Sustainable Development Goals using machine learning. *Nat. Mach. Intell.* **2**,
374 559–565 (2020).
- 375 9. Nunez-Mir, G. C., Iannone, B. V., Curtis, K. & Fei, S. Evaluating the evolution of forest restoration
376 research in a changing world: a “big literature” review. *New For.* **46**, 669–682 (2015).
- 377 10. Westgate, M. J. *et al.* Software support for environmental evidence synthesis. *Nat. Ecol. Evol.* **2**,
378 588–590 (2018).
- 379 11. Lamb, W. F., Creutzig, F., Callaghan, M. W. & Minx, J. C. Learning about urban climate solutions
380 from case studies. *Nat. Clim. Change* **9**, 279–287 (2019).
- 381 12. Cohen, A. M. An Effective General Purpose Approach for Automated Biomedical Document
382 Classification. *AMIA. Annu. Symp. Proc.* **2006**, 161–165 (2006).
- 383 13. Marshall, I. J., Kuiper, J., Banner, E. & Wallace, B. C. Automating Biomedical Evidence Synthesis:
384 RobotReviewer. *Proc. Conf. Assoc. Comput. Linguist. Meet.* **2017**, 7–12 (2017).
- 385 14. Baclic, O. *et al.* Challenges and opportunities for public health made possible by advances in natural
386 language processing. *Can. Commun. Dis. Rep.* **46**, 161–168 (2020).
- 387 15. Schleussner, C.-F. & Fyson, C. L. Scenarios science needed in UNFCCC periodic review | Nature
388 Climate Change. *Nat. Clim. Change* **10**, (2020).
- 389 16. Fankhauser, S. Adaptation to Climate Change. *Annu. Rev. Resour. Econ.* **9**, 209–230 (2017).
- 390 17. Bedsworth, L. W. & Hanak, E. Adaptation to Climate Change. *J. Am. Plann. Assoc.* **76**, 477–495
391 (2010).
- 392 18. IPCC. *Managing the Risks of Extreme Events and Disasters to Advance Climate Change Adaptation*.
393 (Cambridge University Press, 2012).
- 394 19. Hallegatte, S. & Mach, K. J. Make climate-change assessments more relevant. *Nat. News* **534**, 613
395 (2016).
- 396 20. Conway, D. *et al.* The need for bottom-up assessments of climate risks and adaptation in climate-
397 sensitive regions. *Nat. Clim. Change* **9**, 503–511 (2019).

- 398 21. Hansen, G. & Stone, D. Assessing the observed impact of anthropogenic climate change. *Nat. Clim.*
399 *Change* **6**, 532–537 (2016).
- 400 22. Knutson, T. R., Zeng, F. & Wittenberg, A. T. Multimodel Assessment of Regional Surface
401 Temperature Trends: CMIP3 and CMIP5 Twentieth-Century Simulations. *J. Clim.* **26**, 8709–8743
402 (2013).
- 403 23. Knutson, T. R. & Zeng, F. Model Assessment of Observed Precipitation Trends over Land Regions:
404 Detectable Human Influences and Possible Low Bias in Model Trends. *J. Clim.* **31**, 4617–4637
405 (2018).
- 406 24. Nerem, R. S. *et al.* Climate-change-driven accelerated sea-level rise detected in the altimeter era.
407 *Proc. Natl. Acad. Sci. U. S. A.* **115**, 2022–2025 (2018).
- 408 25. Gudmundsson, L., Leonard, M., Do, H. X., Westra, S. & Seneviratne, S. I. Observed Trends in
409 Global Indicators of Mean and Extreme Streamflow. *Geophys. Res. Lett.* **46**, 756–766 (2019).
- 410 26. Padrón, R. S. *et al.* Observed changes in dry-season water availability attributed to human-induced
411 climate change. *Nat. Geosci.* **13**, 477–481 (2020).
- 412 27. Devlin, J., Chang, M.-W., Lee, K. & Toutanova, K. BERT: Pre-training of Deep Bidirectional
413 Transformers for Language Understanding. *ArXiv181004805 Cs* (2019).
- 414 28. Sanh, V., Debut, L., Chaumond, J. & Wolf, T. DistilBERT, a distilled version of BERT: smaller,
415 faster, cheaper and lighter. *ArXiv191001108 Cs* (2020).
- 416 29. Halterman, A. Mordecai: Full Text Geoparsing and Event Geocoding. *J. Open Source Softw.* **2**, 91
417 (2017).
- 418 30. Lane, J. E., Kruuk, L. E. B., Charmantier, A., Murie, J. O. & Dobson, F. S. Delayed phenology and
419 reduced fitness associated with climate change in a wild hibernator. *Nature* **489**, 554–557 (2012).
- 420 31. Zhang, Y. Q., Yu, C. H. & Bao, J. Z. Acute effect of daily mean temperature on ischemic heart
421 disease mortality: a multivariable meta-analysis from 12 counties across Hubei Province, China.
422 *Zhonghua Yu Fang Yi Xue Za Zhi* **50**, 990–995 (2016).

- 423 32. Barry, A. A. *et al.* West Africa climate extremes and climate change indices. *Int. J. Climatol.* **38**,
424 e921–e938 (2018).
- 425 33. Hegerl, G. C. *et al.* Good Practice Guidance Paper on Detection and Attribution Related to
426 Anthropogenic Climate Change. in *Meeting Report of the Intergovernmental Panel on Climate*
427 *Change Expert Meeting on Detection and Attribution of Anthropogenic Climate Change* (eds.
428 Stocker, T. F. *et al.*) (IPCC Working Group I Technical Support Unit, University of Bern, 2010).
- 429 34. Rosenzweig, C. *et al.* Assessment of observed changes and responses in natural and managed
430 systems. in *Climate Change 2007: Impacts, Adaptation and Vulnerability. Contribution of Working*
431 *Group II to the Fourth Assessment Report of the Intergovernmental Panel on Climate Change* 79–
432 131 (Cambridge University Press).
- 433 35. Rosenzweig, C. *et al.* Attributing physical and biological impacts to anthropogenic climate change.
434 *Nature* **453**, 353–357 (2008).
- 435 36. Center for International Earth Science Information Network - CIESIN - Columbia University.
436 Gridded Population of the World, Version 4 (GPWv4): Population Density, Revision 11. (2018).
- 437 37. Frank, D. *et al.* Effects of climate extremes on the terrestrial carbon cycle: concepts, processes and
438 potential future impacts. *Glob. Change Biol.* **21**, 2861–2880 (2015).
- 439 38. Schleussner, C.-F. *et al.* 1.5°C Hotspots: Climate Hazards, Vulnerabilities, and Impacts. *Annu. Rev.*
440 *Environ. Resour.* **43**, 135–163 (2018).
- 441 39. Peng, R. D. Reproducible Research in Computational Science. *Science* **334**, 1226–1227 (2011).
- 442 40. Müller-Hansen, F., Callaghan, M. W. & Minx, J. C. Text as big data: Develop codes of practice for
443 rigorous computational text analysis in energy social science. *Energy Res. Soc. Sci.* **70**, 101691
444 (2020).
- 445 41. Shepherd, T. G. Storyline approach to the construction of regional climate change information. *Proc.*
446 *R. Soc. Math. Phys. Eng. Sci.* **475**, 20190013 (2019).
- 447 42. Rosenzweig, C. & Neofotis, P. Detection and attribution of anthropogenic climate change impacts.
448 *WIREs Clim. Change* **4**, 121–150 (2013).

- 449 43. Mengel, M., Treu, S., Lange, S. & Frieler, K. ATTRICI 1.0 - counterfactual climate for impact
450 attribution. *Geosci. Model Dev. Discuss.* 1–26 (2020) doi:<https://doi.org/10.5194/gmd-2020-145>.
- 451 44. Gudmundsson, L. *et al.* Globally observed trends in mean and extreme river flow attributed to
452 climate change. *Science* **371**, 1159–1162 (2021).
- 453 45. Diffenbaugh, N. S. Verification of extreme event attribution: Using out-of-sample observations to
454 assess changes in probabilities of unprecedented events. *Sci. Adv.* **6**, eaay2368 (2020).
- 455 46. Herring, S. C., Christidis, N., Hoell, A., Hoerling, M. P. & Stott, P. A. Explaining Extreme Events of
456 2019 from a Climate Perspective. *Bull. Am. Meteorol. Soc.* **102**, S1–S116 (2021).
- 457 47. *Cochrane Handbook for Systematic Reviews of Interventions.* (John Wiley & Sons, 2019).
- 458
459
460
461

462 Methods

463 Outline

464 An overview of each of the steps taken in this study is given in Extended Fig 1. These are outlined briefly
465 here and explained in detail in the following sections. Over 600,000 documents were retrieved from
466 bibliographic databases using a query. 2,373 of these documents were screened for relevance and coded
467 for impact type and driver by human reviewers. The implicit inclusion and coding decisions for a further
468 351 documents were extracted from Tables 18.5-18.9 in the contribution of Working Group II to the Fifth
469 Assessment Report of the IPCC¹. Machine learning classifiers were trained to predict relevance of
470 documents using the titles and abstracts, and evaluated using nested cross-validation. The best performing
471 classifier was then fit with all labelled documents using bootstrapping to make predictions with
472 confidence intervals for the relevance of the remaining documents. Those documents predicted to be
473 irrelevant were discarded, as were documents labelled by reviewers as irrelevant. Multilabel classifiers
474 were then trained using the remaining labelled relevant documents, and assessed in a similar fashion using
475 cross-validation. Predictions for impact type and driver were then made for the remaining unlabelled
476 documents. Geographical entities were extracted from the included studies using a geoparser, and each
477 entity was matched to the set of 2.5 degree grid cells overlapping it. Observed trends in precipitation and
478 temperature were collected for 2.5 and 5 degree grid cells and compared with climate models to assess
479 whether observed trends were detectable (i.e., unusual compared with natural variability, and in the same
480 direction as simulated by historical forcing climate model simulations) and at least partially attributable to
481 human influence on the climate, as discussed below. Finally, documents predicted to be driven by
482 temperature or precipitation were extracted from the database of studies and merged with the grid cell
483 attribution datasets so that each document could be characterised by the presence of human-attributable
484 climate trends in the grid cells it referred to, and each grid cell could be characterised by the number of
485 studies referring to it.

486
487

488 Search, screening and coding

489 Search Strategy

490 Potentially relevant documents were assembled by developing a query to search bibliographic databases.
491 To validate the query, we tested this against a set of records known to be relevant. Tables 18.5-18.9 in the
492 contribution of Working Group II to the Fifth Assessment Report of the IPCC¹ (AR5 WGII) contain the
493 studies considered in their assessment of the observed impacts of climate change. After extracting these
494 references, we built a query that would return all of the references in the tables that specifically referred to
495 the role of climate change (rather than of counterfactual explanations for impacts). The query is
496 reproduced in the Supplementary Information (in the format for Web of Science - the same query was
497 used for Scopus) and is made up of three lists of keywords linked with boolean ANDs. The first set of
498 keywords refer to climate and climate variables, the second to impacts, and the third to observations and
499 attribution.

500

501 The query was performed on Scopus and the following citation indices from the Web of Science Core
502 Collection:

- 503 ● Science Citation Index Expanded (SCI-EXPANDED) --1900-present
- 504 ● Social Sciences Citation Index (SSCI) --1900-present
- 505 ● Arts & Humanities Citation Index (A&HCI) --1975-present
- 506 ● Conference Proceedings Citation Index- Science (CPCI-S) --1990-present
- 507 ● Conference Proceedings Citation Index- Social Science & Humanities (CPCI-SSH) --1990-present
- 508 ● Emerging Sources Citation Index (ESCI) --2015-present

509 The queries were updated on October 19 2020: Web of Science returned 411,194 documents, while
510 Scopus returned 476,778 documents. The total number of records after deduplication through fuzzy title
511 and publication year matching using trigram similarity was 601,667. The queries were imported into a
512 database and deduplicated using the NACSOS review platform⁴⁸.

513 Inclusion and exclusion criteria

514 We take a broad definition of climate impacts to include all studies relevant to understanding the observed
515 impacts of climate change. This includes

- 516 ● Studies which explicitly link impacts to climate change (8% of coded studies)
- 517 ● Studies which link impacts to trends in climate drivers like temperature or precipitation (42% of
518 coded studies)
- 519 ● Studies which link impacts to extreme climate events (6% of coded studies)
- 520 ● Studies which link impacts to variation in climate drivers (39% of coded studies)
- 521 ● Studies which document regional or local climate trends (11% of coded studies)

522

523 Documents which only provide evidence of likely future impacts of climate change were excluded.

524

525 With this broad definition of climate impacts evidence, we do not claim that each study is in and of itself
526 evidence of the impacts of climate change. Rather, taken together, and in the context of observations and
527 climate models, this collection of included studies constitutes the evidence base necessary for
528 understanding climate impacts.

529

530 Coding impacts and drivers

531 Where documents were selected for inclusion, reviewers coded the attribution category, the climate
532 impacts and the drivers (where appropriate) for each paper. Impacts and their drivers were chosen from a
533 selection of 75 specific categories, which were aggregated according to the hierarchy of categories
534 included in the supplementary file category_aggregation.csv. 93% of included studies coded impacts in
535 one or more of the 5 broad impact categories used by IPCC AR5:

- 536 ● Mountains, snow and ice (11.42% of included studies)
- 537 ● Rivers, lakes and soil moisture (21.27% of included studies)
- 538 ● Terrestrial ecosystems (33.13% of included studies)
- 539 ● Coastal and marine ecosystems (13.21% of included studies)
- 540 ● Human and managed systems (21.42% of included studies)

541 Remaining studies documented only trends in climate variables without reference to any of these systems.
542

543 Screening and Coding

544
545 A total of 2,373 documents were screened by members of the author team using the NACSOS platform⁴⁸,
546 of which 1,125 were included as relevant and coded for impacts and drivers. The median number of
547 documents coded per user was 133, and the mean was 173.

548
549 In addition, documents extracted from the tables 18.5-18.9 in AR5 WGII were automatically labelled as
550 relevant and tagged with the broad impact categories corresponding to the table in which they were found.

551
552 In order to mitigate a highly unbalanced sample (few relevant documents among many irrelevant
553 documents), and to make best use of reviewing resources, some documents were selected for screening
554 using an adapted active learning pipeline. With active learning, a classifier (see following section for
555 details) is trained using existing screening decisions to predict the relevance of documents yet to be
556 reviewed. Usually, reviewers screen subsequent documents in decreasing order of predicted relevance and
557 the classifier is periodically updated with the new data that has been generated. Given that our goal was to
558 not to screen all relevant documents but to generate useful labels efficiently, we created samples with
559 relevance predictions greater than 0.2, 0.3 and 0.4, in order to exclude documents with a low likelihood of
560 being relevant. Documents were first screened by a small group of reviewers who developed the
561 categorisation scheme for impacts and drivers. A subsequent set of documents was screened by all
562 reviewers, and differences in coding were discussed and alterations recorded. Reviewers were then split
563 into teams corresponding with the AR5 impact categories according to expertise, and screened documents
564 predicted to be rather relevant (>0.33) to the given category. Each team screened a sample of documents
565 and discussed differences in screening and coding decisions. Teams reached average Cohen's Kappa
566 scores between 0.66, indicating substantial agreement, and 1.0, indicating full agreement⁴⁹. After this
567 initial round of double coding, reviewers proceeded to screen documents individually. Additional
568 documents were selected for screening using keyword searches ([https://github.com/mcallaghan/regional-
569 impacts-map/blob/master/literature_identification/category_keywords.ipynb](https://github.com/mcallaghan/regional-impacts-map/blob/master/literature_identification/category_keywords.ipynb)) to identify documents from
570 infrequently appearing subcategories.

571
572 Because the documents selected using the methods described above are unlikely to be representative of
573 the full set of documents returned by the query, we also screened 732 documents drawn at random which
574 we used for validation.

575

576

577

578 Machine learning classifiers for inclusion, impact type and drivers

579 We first trained a binary classifier to predict the inclusion/exclusion decision given by reviewers. We use
580 a nested cross-validation procedure (Extended Fig. 2) to optimize parameter settings and evaluate the

581 performance of a support vector machine (SVM) classifier⁵⁰ as well as a pre-trained DistilBERT model
582 fine-tuned with our labelled dataset²⁸. Support vector machines have a long history of applications in
583 evidence synthesis¹², while the BERT²⁷ (Bidirectional Encoder Representations from Transformers)
584 model recently achieved state of the art results in a variety of natural language processing challenges, and
585 has begun to be used in evidence synthesis pipelines⁸. However, large language models like BERT can
586 have significant climate impacts⁵¹; motivating our decision to use the lighter and faster DistilBERT,
587 which retains “97% of its language understanding”²⁸, with greatly reduced computational resource usage.
588

589 In our nested cross-validation procedure, we first separate those documents which were drawn at random
590 from the population of documents identified by the query from the remaining unrepresentative
591 documents. Only randomly selected documents are used in validation and test sets, in order to ensure that
592 the estimation of the performance of the classifier on the whole dataset is not biased. In the outer fold of
593 the cross-validation loop, a separate test set is drawn from the randomly selected documents for each fold,
594 k, and all other documents are assigned to the test set. The inner CV loop draws k inner validation sets
595 from the remaining random documents in the training set, and allocates all other documents in the training
596 set to an inner training set. The inner loop is used to optimise hyperparameters for each model using grid
597 search: a model is initialised with each combination of hyperparameters and fit on each inner training set
598 and evaluated on each inner validation set. The combination of hyperparameters with the best mean F1
599 score across inner folds is selected as the best model. This model is fit with the training data from the
600 outer CV and evaluated with the test data. The outer CV thus returns k scores for each metric, which we
601 report below. We note that our cross-validation approach, while transparent, robust and thorough, is
602 computationally expensive - and that alternative procedures such as random search may provide similar
603 results at lower computational cost, or minor improvements at the same cost⁵². In principle, additional
604 improvements to the model may also be generated through additional pre-training⁵³ using the unlabelled
605 corpus of climate-relevant abstracts. Pre-training BERT-like models on climate science corpora remains
606 an area for future investigation.
607

608 We evaluated our binary inclusion/exclusion classifiers with 5 inner and outer folds. DistilBERT clearly
609 outperformed SVM across all metrics, achieving an average F1 score of 0.71, and an average ROC AUC
610 score of 0.92 (Extended Fig. 3). A final DistilBERT model configuration was chosen using the same
611 procedure on the outer folds. Each combination of parameter settings was tested on each outer fold, and
612 the combination of parameter settings with the highest mean F1 score was selected.
613

614 This final model was used to predict the relevance of all remaining documents. To create a confidence
615 interval for each prediction, 5 versions of the final model were trained on 5 folds of the data. Upper and
616 lower estimates for each document are given by the mean plus or minus one standard deviation. All
617 documents where the lower estimate was below 0.5 were excluded from the study.
618

619 We then trained multilabel classifiers to predict the impact category and the driver category of included
620 documents. Classifiers parameters were optimised and classifiers evaluated with the same nested cross-
621 validation method, using only those labelled documents which were included. Because documents
622 selected for screening using the active learning process are broadly representative of the documents to
623 which the multilabel classifiers are applied, all documents selected in this manner are also used for
624 validation. Due to the lower number of documents, and lower number of documents drawn from a random

625 sample in this set, we used a smaller k value of 3 for cross-validation. We treat each class equally and
626 optimise using the macro F1 score. For the prediction of impact categories, DistilBERT outperforms
627 SVM, achieving a macro-averaged F1 score of 0.84 and a macro-averaged ROC AUC score of 0.95
628 (Extended Fig. 4.). For classification of climate drivers, we optimise for the macro-averaged F1 score for
629 the categories temperature and precipitation. DistilBERT outperforms SVM, achieving an average F1
630 score of 0.79 and an average ROC AUC score of 0.86. Where no individual class has a prediction larger
631 than 0.5, documents are classes as “Other systems”.

632 Detection and Attribution

633 To put our database of impact studies in context, we match studies with grid cell level detection and
634 attribution of temperature and precipitation trends to human influence on the climate.

635 Updating attribution of temperature and precipitation trends

636 We followed a previously published methodology^{22,23} used to attribute observed temperature and
637 precipitation trends to human influence around the globe, at the level of typical climate model grid cells
638 (5 degree grid boxes for temperature and 2.5 degree grid boxes for precipitation). The different
639 resolutions are based on the available observed datasets, which we did not regrid for our project. The
640 method relies on a comparison of gridbox-scale trends in observational datasets for temperature
641 (HadCRUT4 version 4.6⁵⁴) and precipitation (GPCC v2018, obtainable from
642 <https://psl.noaa.gov/data/gridded/data.gpcc.html>), with those produced in climate model runs from
643 CMIP6⁵⁵. The CMIP6 runs simulate climate changes over the historical period under the influence of
644 either all forcings (i.e., both natural and anthropogenic, referred to as “ALL”) or natural forcings only
645 (referred to as “NAT”).

646
647 We analysed the outputs of these simulations from 10 CMIP6 models, namely MIROC6, IPSL-CM6A-
648 LR, CanESM5, HadGEM3-GC31-LL, CNRM-CM6-1, GFDL-ESM4, CCESS-ESM1-5, BCC-CSM2-
649 MR, NorESM2-LM and CESM2. The model selection was based on the availability of ALL, NAT as
650 well as “piControl” runs (simulating internal climate variations in the absence of external forcings, apart
651 from a constant solar forcing). The analysis provides a test of the ability of the corresponding ALL
652 simulations to reproduce the regional trends in annual mean temperature and precipitation against
653 observational data⁵⁶. For some models the ALL simulations were not available after 2014, in which case
654 we combined them with the first few years of the ssp585 simulations of future climate conditions in order
655 to match the length of the observational data.

656
657 Linear trends over the 1951-2018 (for temperature) and 1951-2016 periods (for precipitation) were
658 computed over each grid cell with adequate data for each observational dataset, following the criteria of
659 ref. 7 and 8 (see Extended Figures 6a&b). For temperature we computed a linear trend for each ensemble
660 member of the HadCRUT4 dataset, from which observed trend distributions were derived. Precipitation
661 trends were not computed over grid cells where less than 20% of data was available for the first or last
662 10% of the observed time series or where the entire time series had less than 70% of data available. For
663 temperature, we divide the trend period into five roughly equal periods and require that each period has at

664 least 20% temporal coverage for annual means. We consider an annual mean as available if at least 40%
665 of the months are available for the year.

666
667 To be compared with the observational data, for each model the data from both the ALL and NAT runs
668 were first re-gridded onto the observational grids ($5^\circ \times 5^\circ$ for temperature and $2.5^\circ \times 2.5^\circ$ for
669 precipitation), excluding times and grid locations where observed data were missing, before linear trends
670 were computed over each grid cell in which adequate temporal coverage was available (see Extended
671 Figures 6c&d). For each model, we then assessed the potential effect of internal variability by computing
672 trends of the length being investigated in 50 random samples of the corresponding piControl runs from
673 each model. The model control runs had beforehand been corrected for any long-term drift, and the
674 anomaly series adjusted by a factor to ensure consistency of low-frequency variability between model
675 control runs and estimated internal variability from observations (further discussed below). We then
676 combined the resulting trend distributions from the piControl runs with the trends computed in the
677 ensemble mean of ALL and NAT runs. Following previous studies^{22,23}, the final trend distribution for
678 temperature was based on an aggregate distribution of all constructed model trend distributions (and thus
679 included the spread of different model ensemble means) whereas for precipitation, an average distribution
680 of model trends across the ensemble was used (i.e., the distribution had the average characteristics of the
681 10 CMIP6 models).

682
683 Attribution categories were assigned to grid cells (Extended Fig. 6 e,f) based on where their observed
684 trend (or trend distribution in the case of temperature) lay relative to the final trend distributions derived
685 from the ALL and NAT runs. Over the grid cells where an observed trend was in the same direction
686 (sign) as the mean of the ALL trend distribution and was outside the trend distribution 5th-95th
687 percentile range for the NAT simulations, the observed trend was categorized as -3 (+3), -2 (+2) or -1
688 (+1) depending on whether it was significantly stronger, the same or weaker than the simulated decrease
689 (increase). Categories -3 (+3) and -2 (+2) are defined as decreases (increases) that are detectable and at
690 least partially attributable to anthropogenic forcing, according to our methodology. Categories -1 (+1) are
691 detectable but not attributable. If the observed trend was significantly different from the NAT distribution,
692 but was in the opposite direction to the mean of the All-Forcing distribution, it was categorized as -4
693 (observed decrease, modeled increase) or +4 (observed increase, modeled decrease). All observed trends
694 (or trend distributions, in the case of temperature) that intersected with the 5th-95th percentile range of
695 the corresponding trend distributions derived from the NAT runs were categorized as non-detectable, or
696 indistinguishable from natural variability (i.e. category 0). Note that for cases where observed trends or
697 trend distributions had a different sign of the mean trend from that of the trend distribution derived from
698 the ALL runs, but were within the range of the Nat run distribution, the corresponding grid cells were also
699 categorised as non-detectable (category 0).

700
701 Once the grid cells were categorised, in the case of temperature the results were re-gridded to a $2.5^\circ \times$
702 2.5° grid to allow superposition with the categories obtained for precipitation.

703
704 Our analysis requires the internal variability for each grid location and variable to be estimated via model
705 control runs. To compare observed estimated internal variability and trends with those generated by the
706 model control runs, Extended Figs. 7 and 8 show fractional difference maps for estimated internal low-
707 frequency variability (model vs. observed) for each model individually and for the ensemble mean of the

708 modeled variability (the latter being most relevant for our analysis, which is based on combined estimated
709 variability across the models). The observed low-frequency internal variability is estimated by
710 subtracting the multi-model ensemble All-Forcing change from the observations and computing the
711 standard deviation of the annual residuals, after application of a 7-year running mean filter. For models,
712 we use the simulated variability from the various control runs, again smoothed with the 7-year running
713 mean smoother. The averaged internal low-frequency variability comparison plot for precipitation
714 (Extended Fig. 7, top panel) shows reds in most regions indicating that by this measure of internal low-
715 frequency variability, the CMIP6 models actually tend to overestimate observed variability levels. So our
716 detection results for precipitation will tend to be conservative, while conversely, the ability of All-Forcing
717 to be consistent with observations will tend to be liberal, because the modeled spread is relatively wide.
718 However, blue regions are evident in Extended Fig. 7 in some tropical regions, including over Africa and
719 South America, indicating an undersimulation of internal low-frequency variability there. We took the
720 internal variability comparisons vs. observed estimated internal variability in Extended Fig. 7 and
721 adjusted the control run variability and trends by the ratio [Obs. stdev / Model stdev] prior to computing
722 our assessment categories. Results without this variability adjustment (not shown) are broadly similar but
723 show more category -4 (unexplained trends of incorrect sign) over Africa, where internal low-frequency
724 variability appears to be underestimated in models according to this analysis; unadjusted results show
725 slightly less detectable human influence in middle and high latitudes, where internal variability is
726 apparently overestimated in models.

727

728 For surface temperature (extended Fig. 8) the internal variability comparison results vs. observed
729 estimates are similar to those of Knutson et al. 2013 for CMIP3 and CMIP5 with a mixture of results:
730 models tend to simulate more internal variability than the observed estimate in northern mid to high
731 latitudes, typically less than observed over most other ocean regions at lower latitudes, and mixed results
732 over land regions. Whether we include the gridpoint-scale adjustment of simulated internal variability in
733 our detection/attribution analysis or not, the results are similar (unadjusted control run-based assessment
734 not shown). For the assessment of 1951-2018 observed trends (Extended Fig. 6), there are some
735 additional regions with detectable anthropogenic warming compared to Knutson et al. (2013), but that is
736 as expected, since the Knutson et al. analysis only examined trends through 2010. With the termination
737 of the ‘global warming hiatus’ around 2014, the additional recent years have been adding to an ongoing
738 strengthening warming signal and leading to even greater assessed area with detectable anthropogenic
739 warming. In Extended Fig. 6 and elsewhere in the study, we use the adjusted control run results for our
740 assessments for both temperature and precipitation.

741

742 Spatial resolution of studies

743

744 In order to match this data with the finest-scale resolution of our database, we resolved each study to the
745 set of 2.5 degree grid cells contained by the smallest geographical entity extracted from each paper’s title
746 and abstract using the geoparser Mordecai²⁹. For each study, we calculated the proportion of the grid cells
747 that this entity corresponds to in which an attributable trend for each variable can be found. For example,
748 in Extended Figure 9, panels a. and b. show that 20 out of Sudan’s 27 grid cells show an attributable
749 anthropogenic warming trend, so each study referring to Sudan and documenting impacts predicted to be
750 driven by temperature receives a precipitation trend proportion value of 20/27. Such a study would

751 therefore add towards the dark red bars in Fig. 3, which count studies where an attributable temperature
752 trend can be demonstrated for more than 50% of the grid cells the study refers to.

753
754 We also calculate a weighted number of studies for each grid cell by adding 1 divided by the number of
755 grid cells a study refers to each of those grid cells, and repeating this procedure for all identified relevant
756 studies. Extended Figures 9c. and d. show 11 studies which refer to impacts predicted to be driven by
757 temperature trends in Sudan, where Sudan is the smallest geographical entity mentioned. Each gridcell in
758 Sudan therefore receives 11/27 weighted studies. Given that some geographical entities were too small to
759 hold one 2.5 degree grid cell, their longitude-latitude values were interpolated to the nearest grid cell
760 instead and the grouped studies apportioned to that one grid cell. Because 4 additional studies refer to
761 Khartoum, we add 4/1 to the weighted studies value in the grid cell containing Khartoum.
762

763 Data availability statement

764 The results of this study are made available in a public repository⁵⁷

765

766 Code availability statement

767 The code used to produce these results is made available in a public repository⁵⁸

768 Acknowledgements

769 M.C. is supported by a PhD stipend from the Heinrich Böll Stiftung. J.C.M. acknowledges funding from
770 the ERC-2020-SyG GENIE (Grant ID: 951542). S.N. and Q.L. acknowledge funding from the German
771 Federal Ministry of Education and Research (BMBF) and the German Aerospace Center (DLR) via the
772 LAMACLIMA project as part of AXIS, an ERANET initiated by JPI Climate ([http://www.jpi-](http://www.jpi-climate.eu/AXIS/Activities/LAMACLIMA)
773 [climate.eu/AXIS/Activities/LAMACLIMA](http://www.jpi-climate.eu/AXIS/Activities/LAMACLIMA), last access: 26 August 2021, grant no. 01LS1905A), with co-
774 funding from the European Union (grant no. 776608). M.R. acknowledges support by the ERC-SyG
775 USMILE (Grant ID 85518). R.J.B. acknowledges support from the EU Horizon2020 Marie-Curie
776 Fellowship Program H2020-MSCA-IF-2018 (Proposal number 838667 -INTERACTION). We
777 acknowledge the World Climate Research Programme, which, through its Working Group on Coupled
778 Modelling, coordinated and promoted CMIP6. We thank the climate modeling groups for producing and
779 making available their model output, the Earth System Grid Federation (ESGF) for archiving the data and
780 providing access, and the multiple funding agencies who support CMIP6 and ESGF.

781 Author contribution statement

782 M.C. and J.C.M., and C-F.S. designed the research. M.C. developed the coding platform and machine learning
783 pipeline to identify studies, with advice from M.R., M.C, C-F.S., G.H., Q.L., E.T. developed the codebook and
784 coordinated screening and coding. M.C., Q.L., S.N., C-F.S.. conceptualised the link to detection and attribution data.
785 S.N. performed the univariate detection and attribution analysis of temperature and precipitation trends and
786 assessment of internal variability, in consultation with T.R.K, who designed the methodology for these calculations.

787 M.C. and S.N. designed and implemented the matching of studies with detection and attribution data. M.C., C-F.S.,
788 S.N., Q.L, G.H., E.T., M.A., R.J.B., M.H., C.J., K.L., A.L., N.M., I.M., P.P., and B.Y. contributed to screening and
789 coding studies. M.C., C-F.S., J.C.M., Q.L., and S.N. wrote the manuscript with contributions from all authors.

790 Competing interest statement

791 The authors declare no competing interests

792

793 Figure Legends

794

795 **Fig. 1| Results of the machine-assisted literature review.** All results shown are based on our search
796 queries and subsequent classification by the machine learning pipeline. **a.** Growth in the scientific
797 literature relevant to observed climate impacts over the last 30 years (cumulative totals for IPCC
798 assessment periods are highlighted for reference). Inset: numbers of documents considered in the total
799 query and in the IPCC AR5 WGII Tables 18.5-18.9. **b-c.** The estimated number of studies for each impact
800 category and continent in our database (note that uncertainty bars take into account uncertainty over
801 relevance as well as impact category). ES = ecosystem.

802

803 **Fig. 2| Potential attribution of impact studies to regional anthropogenic temperature and**
804 **precipitation trends.** Model-based assessment of the attribution of regional temperature (**a**, timespan
805 1951-2018) and precipitation trends (**b**, timespan 1951-2016) to human influence. Cooling/warming or
806 drying/wetting trends in the regions marked as categories -/+2 and -/+3 are assessed as attributable in part
807 to human influence (see Methods). **c.** Global map of area-weighted studies coloured by the existence of
808 attributable trends (purple for attributable trends in at least one variable, cross-hatched for attributable
809 trends in both variables, grey for no attributable trends) and indicating the localised evidence density
810 (Low: <5 weighted studies, Robust: >5 weighted studies, High: >20 weighted studies). **d.** the proportion
811 of land area and population with each grid cell type, grouped by country income category.

812

813 **Fig 3| A global density map of climate impact evidence.** Map colouring denotes the number of
814 weighted studies per grid cell for all evidence on climate impacts (N=77,785). Bar charts show the
815 number of studies per continent and impact category. Bars are coloured by the climate variable predicted
816 to drive impacts. Colour intensity indicates the percentage of cells a study refers to where a trend in the
817 climate variable can be attributed (partially attributable: >0% of grid cells, mostly attributable: >50% of
818 grid cells).

819

820

821 *CAPTIONS FOR EXTENDED DATA FIGURES*

822

823 **Extended Data Fig. 1 |** A visual representation of the workflow of our machine learning assisted
824 attribution map. Squares represent documents (not to scale), boxes represent the steps taken. Documents
825 are screened by hand, and those labels are used to generate predictions and machine label documents.

826 These machine-labelled documents are matched by location with information from observations and
827 climate models on the detection and attribution of trends in temperature and precipitation.

828 **Extended Data Fig. 2** | Nested cross validation (CV) procedure for the binary relevance classifier.
829 Models are fit using training documents and evaluated on validation/test documents. The inner CV
830 loop is used to search for optimal hyperparameter settings, which are then evaluated on the outer test sets.

831 **Extended Data Fig. 3** | Performance metrics for the binary inclusion/exclusion classifier. Each pair of
832 dots represents the scores for a distinct cross-validation fold. Horizontal lines show the mean score
833 across folds.

834 **Extended Data Fig. 4** | Receiver operating curve area under the curve scores (ROC AUC) and F1 scores
835 for the classification of impact categories. Each pair of dots represents the scores for a distinct cross-
836 validation fold. Horizontal lines show the mean score across folds.

837 **Extended Data Fig. 5** | Receiver operating curves area under the curve scores (ROC AUC)(ROC) and F1
838 scores for the classification of drivers. Each pair of dots represents the scores for a distinct cross-
839 validation fold. Horizontal lines show the mean score across folds.

840 **Extended Data Fig. 6** | Geographical distribution of surface trends. Temperature from 1951 to 2018
841 (left) and precipitation trends from 1951 to 2016 (right) in (a),(b) observations and (c),(d) CMIP6 10-
842 model ensemble mean all-forcing runs. Bottom panels (e),(f) show observations categorised into
843 attribution categories, following refs. 8,7, respectively. Observed cooling/warming or drying/wetting
844 trends that—after accounting for internal climate variability—are inconsistent with the simulated
845 response to natural forcings but consistent with the simulated response to both natural and
846 anthropogenic forcings are indicated by categories -/+2. This is clearest case of changes that are at
847 least partially attributable to anthropogenic forcing, according to the CMIP6 ensemble. Categories -
848 /+1 have detectable observed changes, but are not assessed as attributable to anthropogenic forcing
849 because the observed changes are significantly less than those simulated in the average all-forcing
850 runs. Categories -/+3 have detectable changes and are assessed as at least partly attributable
851 anthropogenic forcing, although the observed changes are inconsistent with the all-forcing runs. That
852 is, they are in the same direction as, but are significantly stronger than, the mean of the all-forcing
853 runs. Categories -/+4 represents cooling/warming or drying/wetting trends that are inconsistent with
854 the simulated response to natural forcings but whose sign is opposite to that of the average simulated
855 all-forcing response; category 0 represents trends that are not distinguishable from natural
856 variability alone. Categories -/+4 and 0 are considered to be examples of non-detectable trends).

857 **Extended Data Fig. 7** | Fractional difference between average CMIP6 modeled low-frequency standard
858 deviation of annual mean precipitation vs observed precipitation. To estimate the internal low-frequency
859 variability for both models and observations, the observed time series were detrended and low-pass
860 filtered with a 7-year running mean filter prior to computing the standard deviations while for the
861 models we used the full available control runs (7-yr running mean filtered) to estimate the internal low-
862 frequency variability for each model. The top panel shows the multi-model ensemble standard deviation
863 comparison while the ten individual panels below it show the comparison for each individual CMIP6
864 model used in the study. The fraction difference was computed as: [(Model st. dev. - Observed st. dev.)
865 / (Observed st. dev.)].

866 **Extended Data Fig. 8** | Difference between average CMIP6 modeled low-frequency standard
867 deviation (°C) of annual mean surface air temperature vs observed surface temperature. To
868 estimate the internal low-frequency variability for both models and observations, the observed
869 time series were detrended and low-pass filtered with a 7-year running mean filter prior to computing

870 the standard deviations while for the models we used the full available control runs (7-year running
871 mean filtered) to estimate the internal low-frequency variability for each model. The top panel shows
872 the multi-model ensemble standard deviation comparison while the ten individual panels below it
873 show the comparison for each individual CMIP6 model used in the study.

874 **Extended Data Fig. 9** | An illustration of the spatial resolution and weighting methodology. Detection
875 and attribution categories for temperature in East Africa; **b.** the number of grid cells of each type in
876 Sudan; **c.** weighted studies for each grid cell in Sudan; **d.** The number of studies referring to each
877 extracted geographical location in Sudan.

878
879

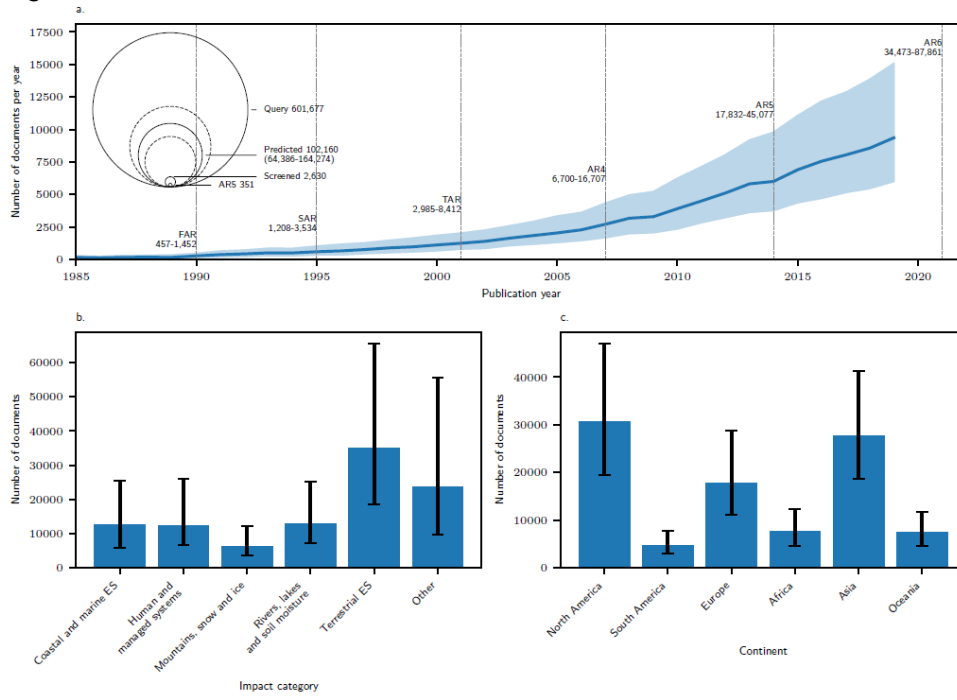
880 References

- 881 48. Callaghan, M., Müller-Hansen, F., Hilaire, J. & Lee, Y. T. *NACSOS: NLP Assisted Classification,*
882 *Synthesis and Online Screening.* (Zenodo, 2020). doi:10.5281/zenodo.4121526.
- 883 49. McHugh, M. L. Interrater reliability: the kappa statistic. *Biochem. Medica* **22**, 276–282 (2012).
- 884 50. Chang, C.-C. & Lin, C.-J. LIBSVM: A library for support vector machines. *ACM Trans. Intell. Syst.*
885 *Technol.* **2**, 27:1-27:27 (2011).
- 886 51. Bender, E. M., Gebru, T., McMillan-Major, A. & Shmitchell, S. On the Dangers of Stochastic
887 Parrots: Can Language Models Be Too Big? 🦜 in *Proceedings of the 2021 ACM Conference*
888 *on Fairness, Accountability, and Transparency* 610–623 (Association for Computing Machinery,
889 2021). doi:10.1145/3442188.3445922.
- 890 52. Bergstra, J. & Bengio, Y. Random Search for Hyper-Parameter Optimization. *J. Mach. Learn. Res.*
891 **13**, 281–305 (2012).
- 892 53. Gururangan, S. *et al.* Don't Stop Pretraining: Adapt Language Models to Domains and Tasks.
893 *ArXiv200410964 Cs* (2020).
- 894 54. Morice, C. P., Kennedy, J. J., Rayner, N. A. & Jones, P. D. Quantifying uncertainties in global and
895 regional temperature change using an ensemble of observational estimates: The HadCRUT4 data set.
896 *Atmospheres* (2012) doi:https://doi.org/10.1029/2011JD017187.
- 897 55. Eyring, V. *et al.* Overview of the Coupled Model Intercomparison Project Phase 6 (CMIP6)
898 experimental design and organization. *Geosci. Model Dev.* **9**, 1937–1958 (2016).

- 899 56. Beusch, L., Gudmundsson, L. & Seneviratne, S. I. Crossbreeding CMIP6 Earth System Models With
900 an Emulator for Regionally Optimized Land Temperature Projections. *Geophys. Res. Lett.* **47**,
901 e2019GL086812 (2020).
- 902 57. Callaghan, M. *et al.* Machine learning-based evidence and attribution mapping of 100,000 climate
903 impact studies - Data. (2021) doi:10.5281/ZENODO.5257271.
- 904 58. Callaghan, M. Machine learning-based evidence and attribution mapping of 100,000 climate impact
905 studies - Code. (2021) doi:10.5281/ZENODO.5327409.
- 906
907

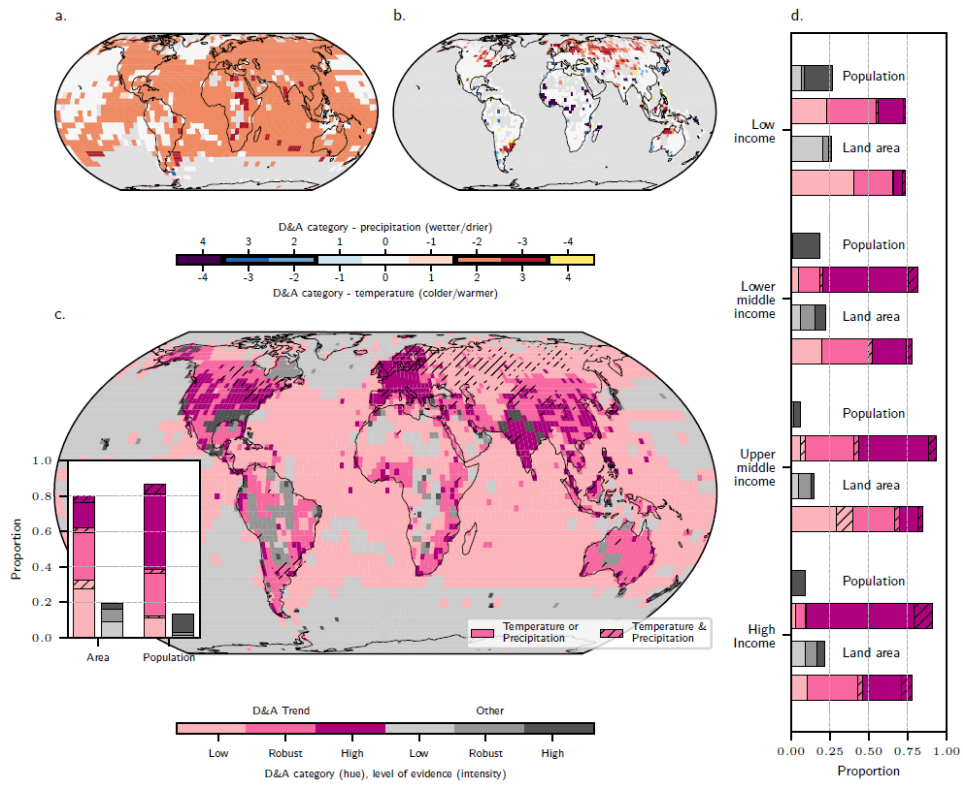
908
909

Figure 1.



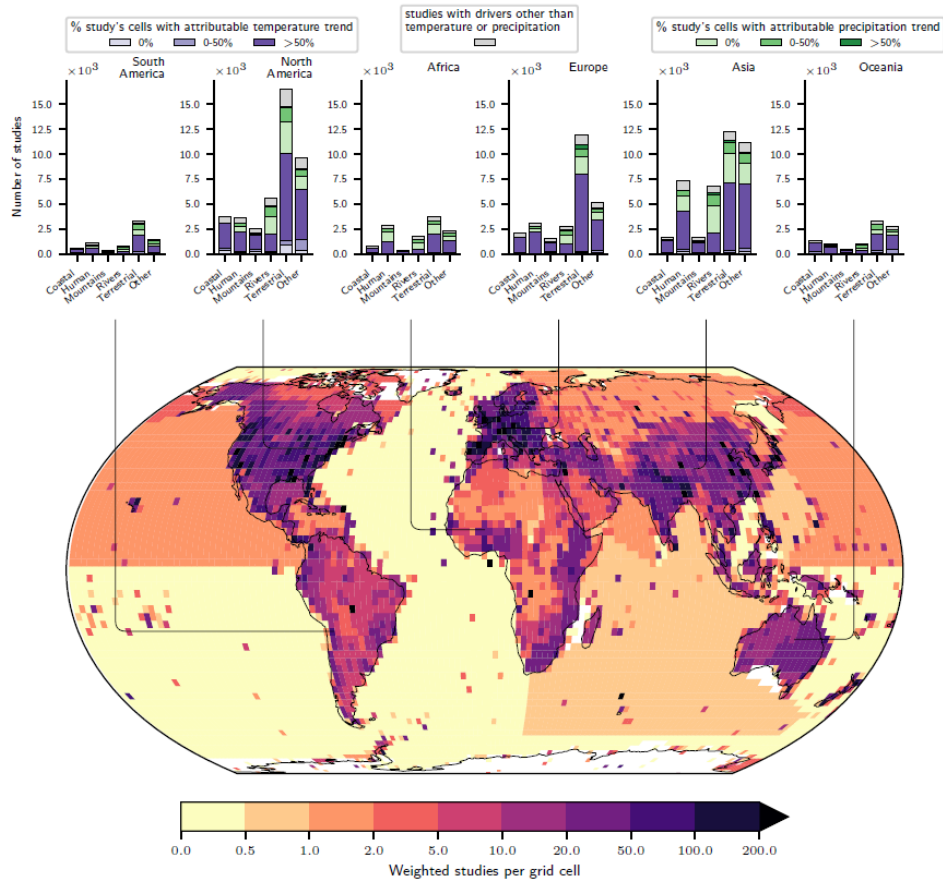
910
911
912

913 Figure 2.



914
915
916

917 Figure 3.



918
919
920
921

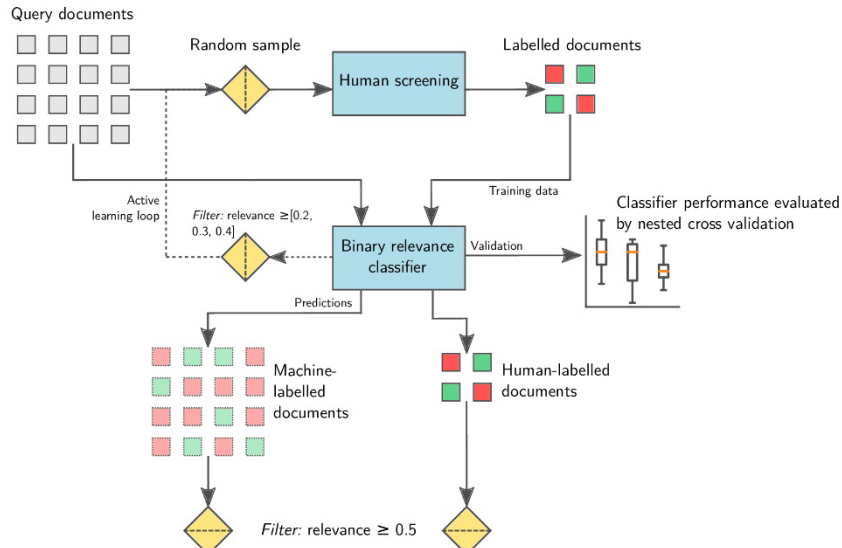
922 Extended Data Figures

923

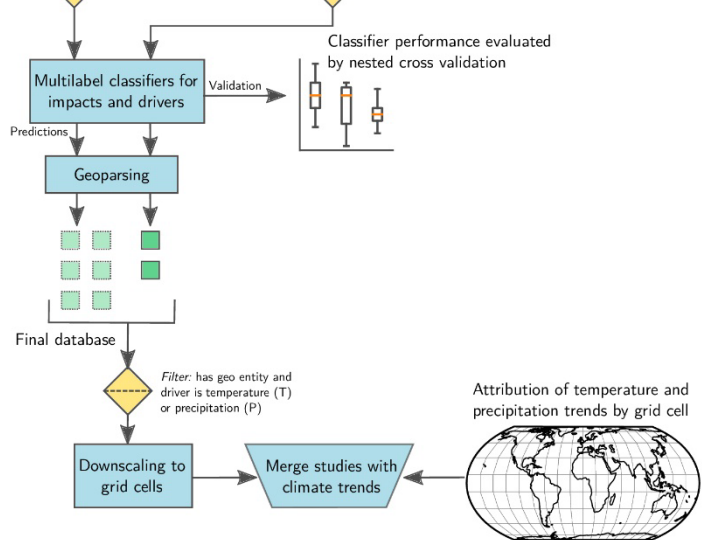
924 Extended Data Fig. 1

925

1. Identification of relevant impacts studies



2. Classification and location of climate impacts

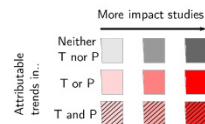


Results:

After merging studies with climate trends, we characterise each study by the proportion of its gridcells which show attributable trends



And each grid cell by the presence of attributable trends in temperature (T) and precipitation (P) and the number of studies on impacts

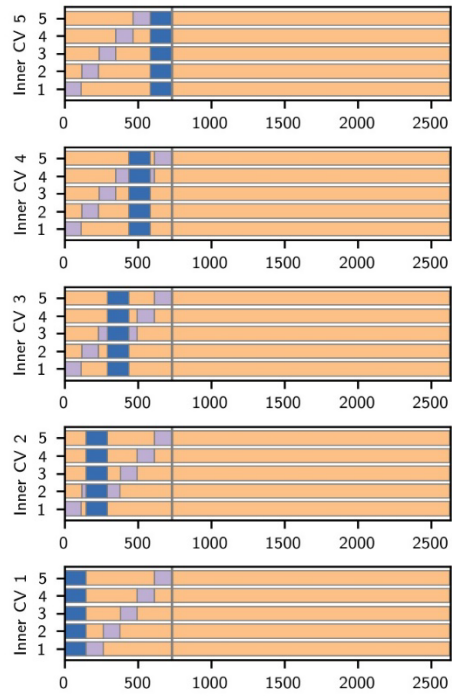
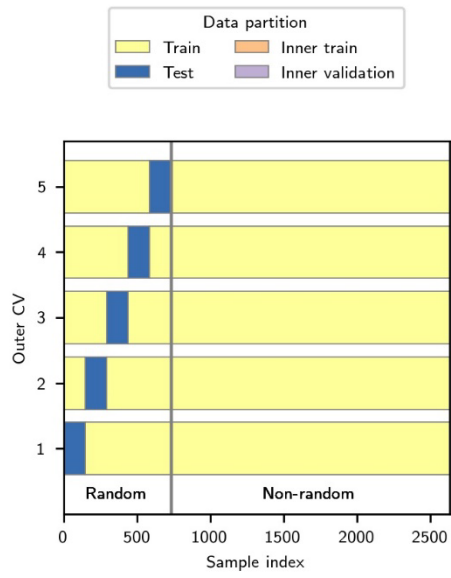


926

927

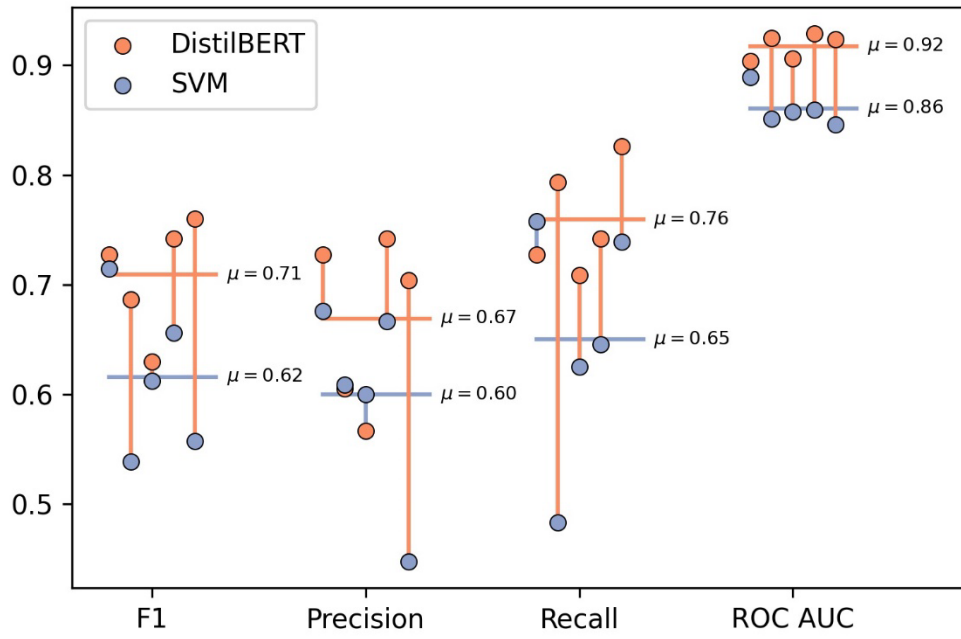
928

929 Extended Data Fig. 2



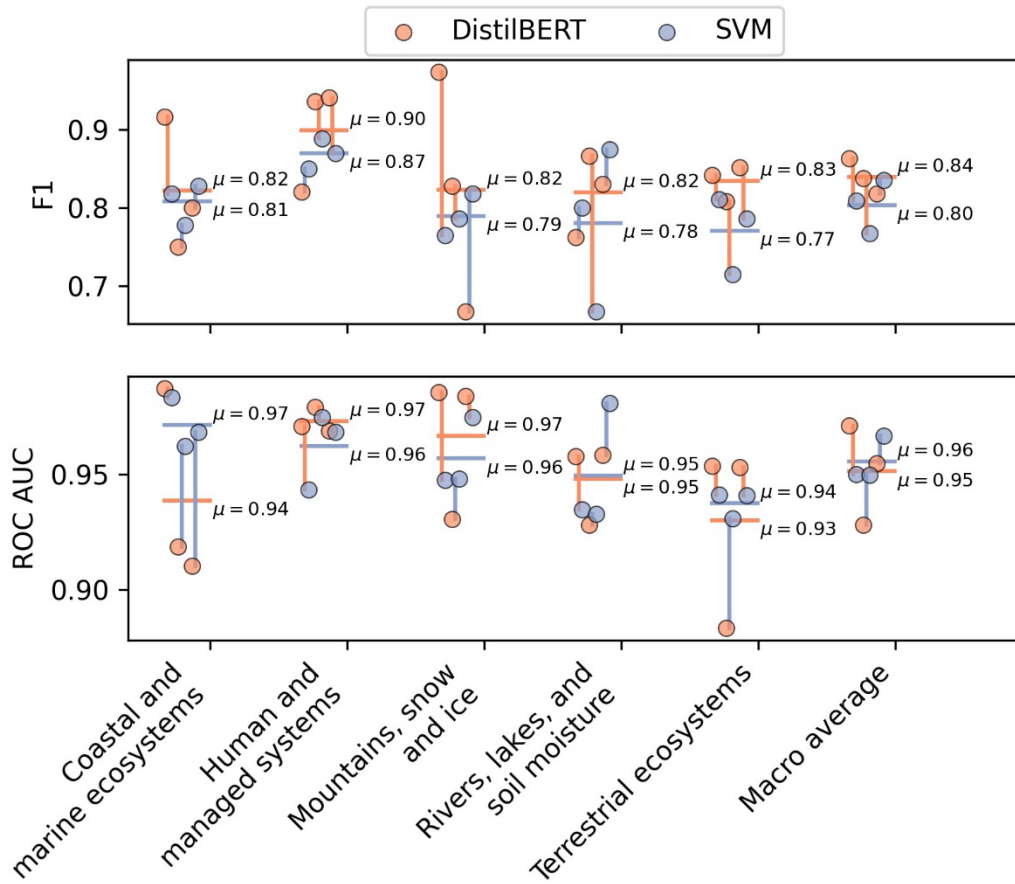
930
931
932

933 Extended Data Fig. 3
934
935
936



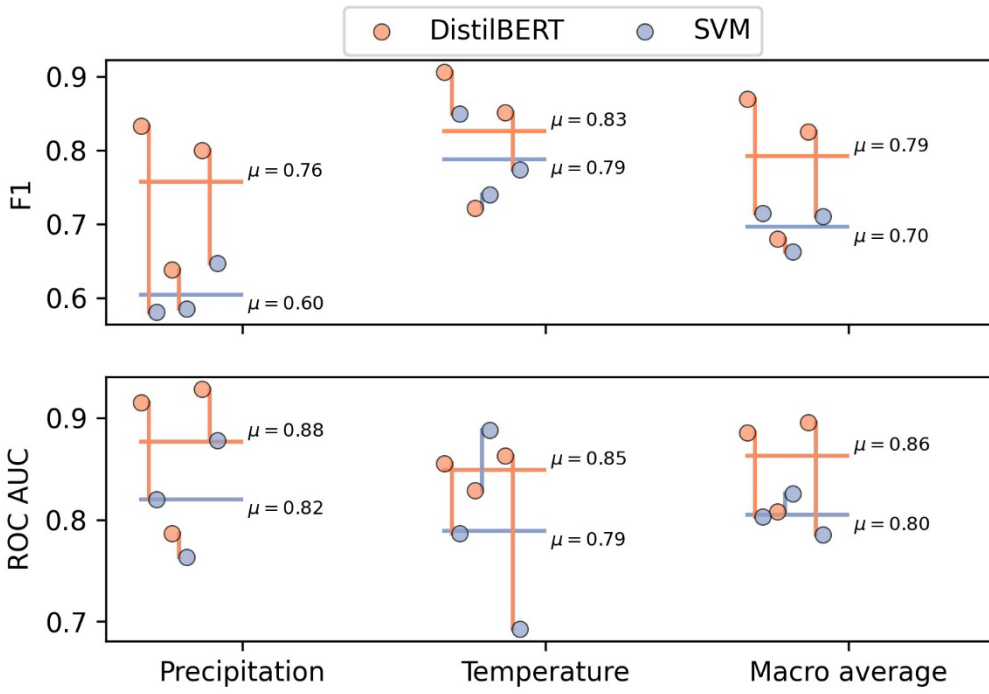
937
938

939 Extended Data Fig. 4



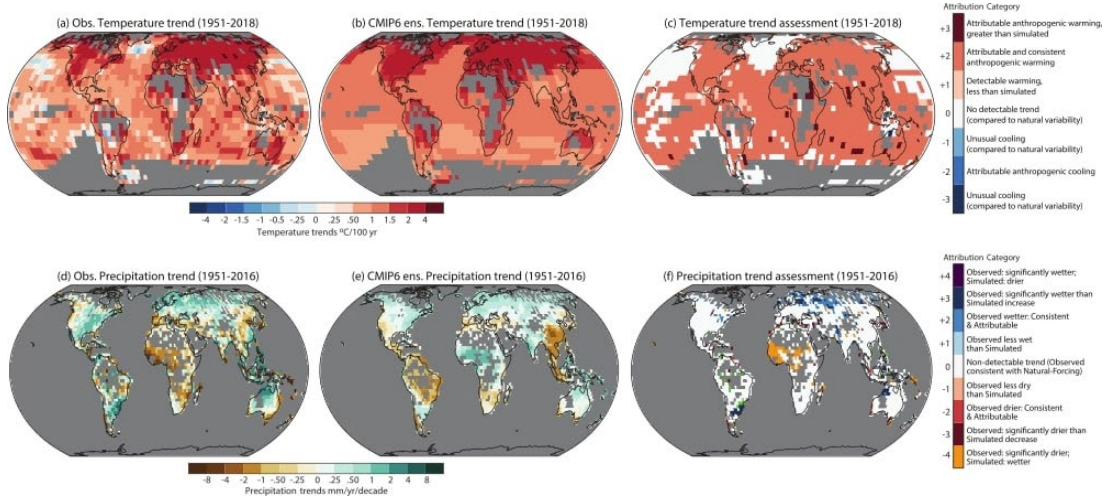
940
941
942

943 Extended Data Fig. 5

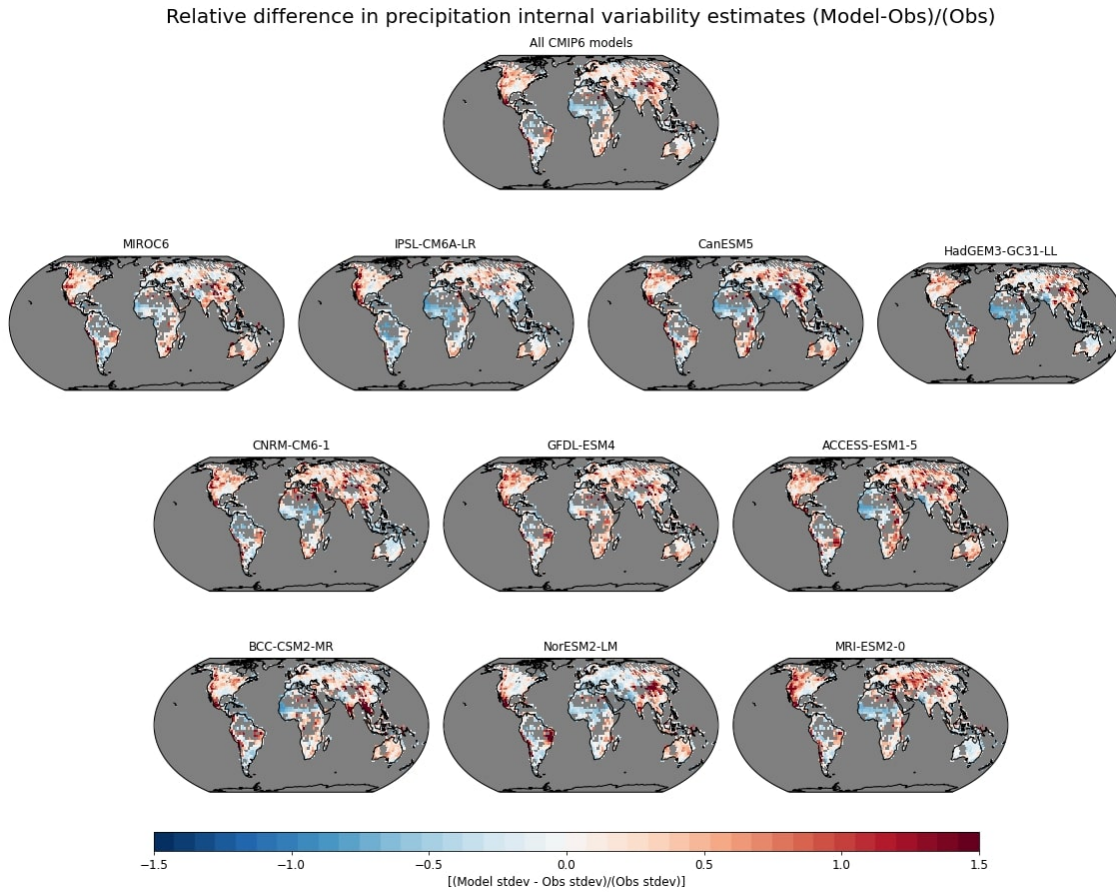


944
945
946
947

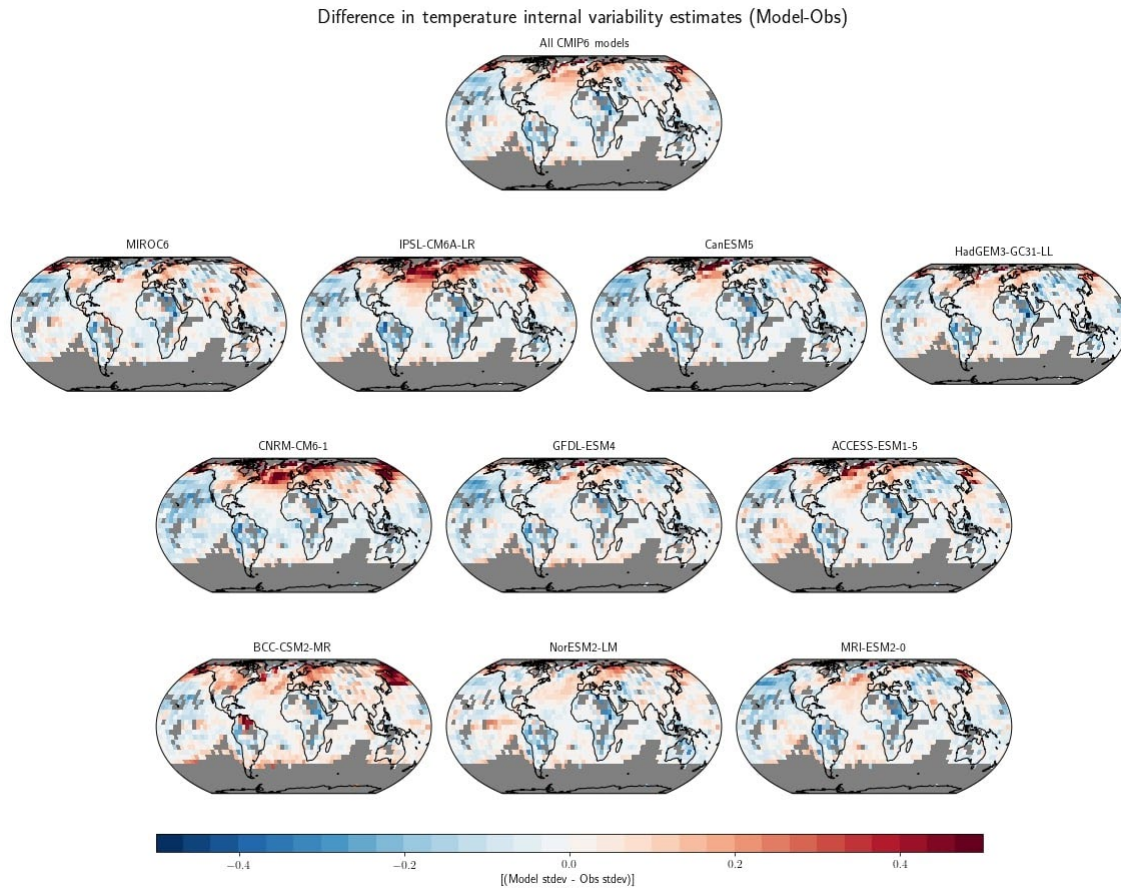
948 Extended Data Fig. 6



949
950
951 Extended Data Fig. 7

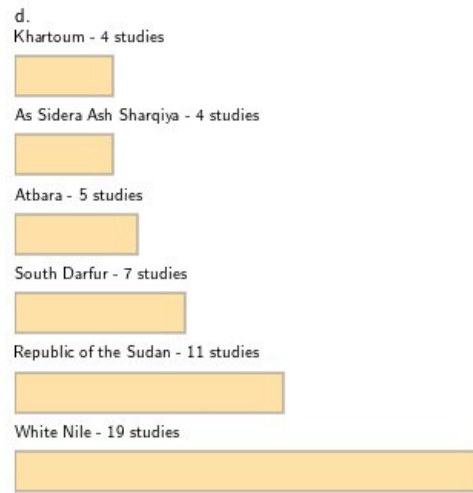
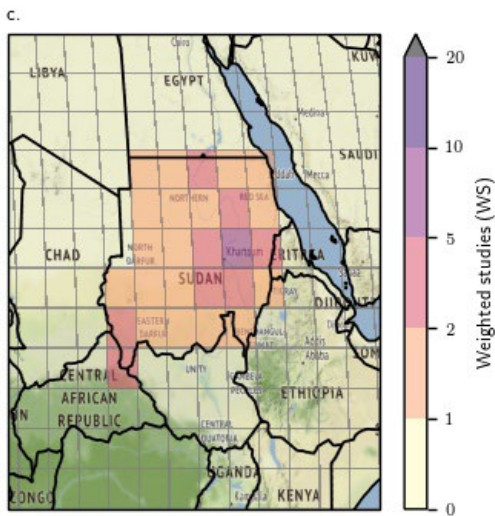
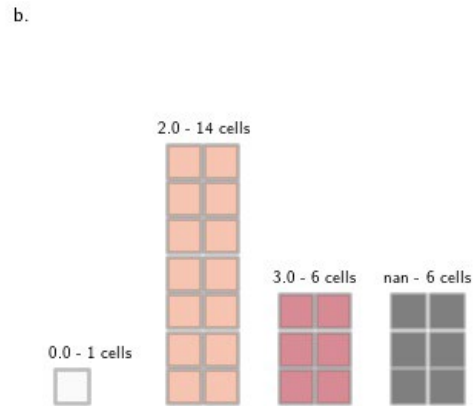
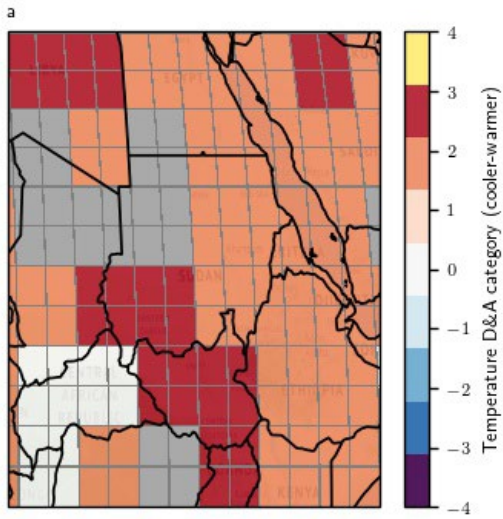


952
953
954
955



957
958
959

960 Extended Data Fig. 9
 961

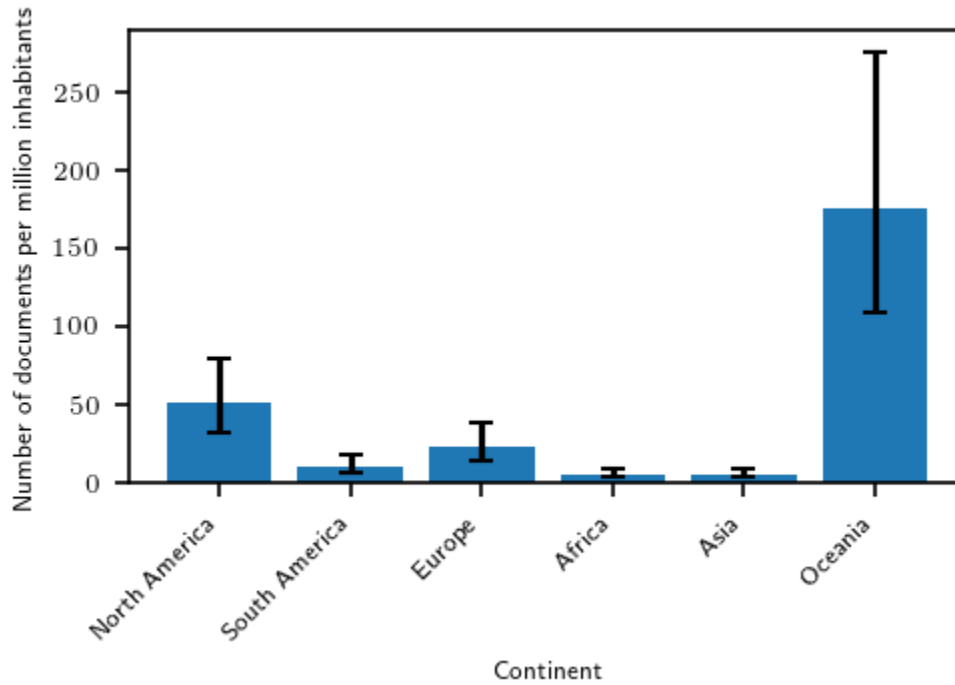


962
 963
 964

965 SUPPLEMENTARY INFORMATION FOR:
966
967 AI based evidence and attribution mapping of 100,000 climate impact studies
968 by Callaghan et al.
969

970 Query

971
972 (TS=("climate model" OR "elevated* temperatur" OR "ocean* warming" OR "saline* intrusion" OR
973 "chang* climat" OR "environment* change" OR "climat* change" OR "climat* warm" OR "warming*
974 climat" OR "climat* varia" OR "global* warming" OR "global* change" OR "greenhouse* effect" OR
975 "snow cover" OR "extreme temperature" OR "cyclone" OR "ocean acidification" OR "anthropogen*" OR
976 "sea* level" OR "precipitation variabil*" OR "precipitation change*" OR "temperature* impact" OR
977 "environmental* variab" OR "weather* pattern" OR "weather* factor*" OR "climat*") OR TS=("change*
978 NEAR/5 cryosphere" OR "increase* NEAR/3 temperatur*"))
979 AND
980 (TS=("migration" OR "impact*" OR "specie*" OR "mortality*" OR "health" OR "disease*" OR
981 "ecosystem*" OR "mass balance" OR "flood*" OR "drought" OR "disease*" OR "adaptation" OR
982 "malaria" OR "fire" OR "water scarcity" OR "water supply" OR "permafrost" OR "biological response"
983 OR "food availability" OR "food security" OR "vegetation dynamic*" OR "cyclone*" OR "yield*" OR
984 "gender" OR "indigenous" OR "conflict" OR "inequality" OR "snow water equival*" OR "surface
985 temp*") OR TS=("glacier* NEAR/3 melt*" OR "glacier* NEAR/3 mass*" OR "erosion* NEAR/5
986 coast*" OR "glacier* NEAR/5 retreat*" OR "rainfall* NEAR/5 reduc*" OR "coral* NEAR/5 stress*" OR
987 "precip* NEAR/5 *crease*" OR "river NEAR/5 flow"))
988 AND
989 (TS=("recent" OR "current" OR "modern" OR "observ*" OR "evidence*" OR "past" OR "local" OR
990 "region*" OR "significant" OR "driver*" OR "driving" OR "respon*" OR "were responsible" OR "was
991 responsible" OR "exhibited" OR "witnessed" OR "attribut*" OR "has increased" OR "has decreased" OR
992 "histor*" OR "correlation" OR "evaluation"))
993
994
995 Supplementary Figures



996
997

998 **Supplementary Fig. 1: The number of papers published in each continent scaled by population.**
999 Bars show the estimated number of papers mentioning a location in each region, uncertainty bars are
1000 generated through bootstrapping (see methods).
1001