



**HAL**  
open science

## Toward an observatory of the evolution of COVID-19 Vaccines Trials through phylomemy reconstruction

Quentin Lobbé, David Chavalarias, Alexandre Delanoë, Gabriel Ferrand,  
Sarah Cohen-Boulakia, Philippe Ravaud, Isabelle Boutron

### ► To cite this version:

Quentin Lobbé, David Chavalarias, Alexandre Delanoë, Gabriel Ferrand, Sarah Cohen-Boulakia, et al.. Toward an observatory of the evolution of COVID-19 Vaccines Trials through phylomemy reconstruction. *Journal of Clinical Epidemiology*, 2022, 149, pp.34-44. 10.1016/j.jclinepi.2022.05.004 . hal-03500847

**HAL Id: hal-03500847**

**<https://hal.science/hal-03500847v1>**

Submitted on 5 Jan 2022

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Toward an observatory of the evolution of COVID-19 Vaccines Trials through phylomemy reconstruction

Quentin Lobbé<sup>a</sup>, David Chavalarias<sup>a,1</sup>, Alexandre Delanoë<sup>a</sup>, Gabriel Ferrand<sup>b,c,d</sup>, Sarah Cohen-Boulakia<sup>e</sup>, Philippe Ravaud<sup>b,c,d</sup>, and Isabelle Boutron<sup>b,c,d</sup>

<sup>a</sup>CNRS, Complex Systems Institute of Paris Île-de-France; <sup>b</sup>Université de Paris, INSERM, INRAE, CNAM, CRESS, F-75004 Paris, France; <sup>c</sup>Centre d'Épidémiologie Clinique, AP-HP, Hôpital Hôtel-Dieu, F-75004 Paris, France; <sup>d</sup>Cochrane France, F-75004 Paris, France; <sup>e</sup>Université Paris-Saclay, France

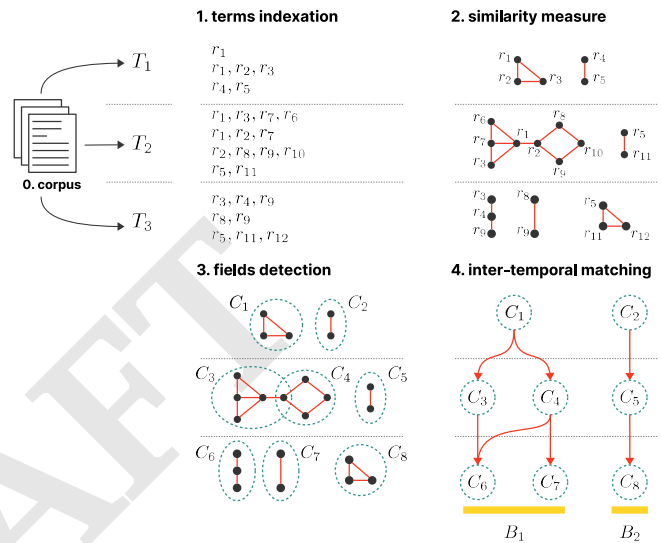
1 This paper aims at reconstructing the evolution of all the available  
 2 COVID-19 vaccines trials extracted from the *COVID-NMA database*  
 3 by applying the *phylomemy reconstruction process*. We visualize  
 4 the textual contents of 1,794 trials descriptions and explore their col-  
 5 lective structure along with their semantic dynamics. We map the  
 6 continuous progress of the main COVID-19 vaccine platforms from  
 7 their early-stage trials in February 2020 to their most recent combi-  
 8 nations driven by the rise of variants of concern, third dose issues  
 9 and heterologous vaccinations. This paper brings insights for the  
 10 global coordination between research teams especially in crisis situ-  
 11 ations such as the COVID-19 pandemic.

COVID-19 | vaccination | phylomemy | knowledge dynamics

1 **Significance statement.** *The COVID-19 pandemic has resulted*  
 2 *in an unprecedented volume of publications that have generated*  
 3 *an information overload for the medical community. One of*  
 4 *today's challenges is to synthesize this overwhelming amount*  
 5 *of information in order to improve coordination between the*  
 6 *different research streams. Our paper thus proposes to apply a*  
 7 *new method for reconstructing the evolution of knowledge and*  
 8 *to visualize the collective structure and semantic dynamics of*  
 9 *1,794 COVID-19 vaccines trials descriptions.*

10 **O**ver the past two years, the ongoing COVID-19 pandemic  
 11 has impacted a wide number of human domains: from  
 12 economy to education, from public health to politics. Among  
 13 others, Science swung early on into action to find both a  
 14 cure and an effective vaccine. This has resulted in an un-  
 15 precedented volume of publications that have generated an  
 16 information overload for the medical community. One of to-  
 17 day's challenges is to synthesize this overwhelming amount  
 18 of information about current COVID-19 research in order to  
 19 improve coordination between the different research streams.  
 20 Our paper thus proposes to address this issue by applying a  
 21 new method for reconstructing the evolution of knowledge.  
 22 We take as a case study the COVID-19 vaccines clinical trials  
 23 from the *COVID-NMA database* and use the *phylomemy re-*  
 24 *construction process* (1). The *COVID-NMA database* stores  
 25 the curated dataset of all the clinical trials available in the set  
 26 of international primary and secondary trial registries\* (2, 3)  
 27 (see [Materials](#)). For the purpose of this study, the *COVID-*  
 28 *NMA database* has been reduced to a pruned corpus called  $\mathcal{D}_{vt}$   
 29 (see [Pre-processing](#)). We then combine the expertise of epi-  
 30 demiologists and *Complex Systems* researchers to interpret  
 31 the resulting visualizations and reveal insights for upcoming  
 32 COVID-19 research.

\* *i.e.*, all trials registered in the International Clinical Trials Registry Platform (ICTRP), Clinicaltrials.gov and the EU clinical trials registry



**Fig. 1.** The four operators of the phylomemy reconstruction process: 1. terms indexation, 2. similarity measures, 3. fields detection, 4. inter-temporal matching

## The phylomemy reconstruction process

The phylomemy reconstruction process (1, 4) combines advanced text-mining methods, scientometrics and methods for the reconstruction of evolving complex networks in order to reconstruct the latent semantic structures of an unstructured – but timestamped – set of textual documents. Applied to a scientific corpus, it results in an inheritance network of research areas covered by all the collected publications. The phylomemy reconstruction process can be described as a combination of four subsequent operators of summarized by the [Figure 1](#):

1. **Terms indexation.** By means of natural language processing (NLP) algorithms and human validations<sup>†</sup>, we first extract from an original corpus of documents ([Figure 1.0](#)) a core vocabulary as a list  $\mathcal{L} = \{r_i \mid i \in \mathcal{I}\}$  of sets  $r_i$  of equivalent expressions called *roots* ([Figure 1.1](#)).

In our case study, the corpus is a set of 1,794 trials descriptions. The *roots* are all the technical and equivalent names (including characteristics variations and any misspelling) given for a same vaccine. For instance, the technical

<sup>†</sup> NLP algorithms and human validations are handled by the free software *Gargantext* (5)

The authors declare no competing interests

expressions “rad5” and “rad26” were aggregated into “gam-covid-vac”<sup>‡</sup>.

The corpus is then sliced into periods of interest  $\mathcal{T}^* = \{T_i\}_{1 \leq i \leq \kappa}$ ,  $T_i \subset \mathcal{T}$  for which roots’ co-occurrences are computed.

In our case study, we consider two weeks periods starting every monday from February 2020 to October 2021 and the output is a series of matrices of roots co-occurrences.

2. **Similarity measure.** Within each period of time and on the basis of its co-occurrences matrix, we estimate the semantic similarity between roots using the *confidence* measure (6). The completion of this task results in a temporal series of graphs of similarity (Figure 1.2).

3. **Fields detection.** For each period, a community detection algorithm – the *frequent item set* method (7) – is applied to detect subsets of densely connected roots within the graphs of similarity. These subsets  $C^T$  are called *fields* (Figure 1.3) and their aggregated root expressions describe consistent research topics that were explored at a given period.

In our case study, the *fields* correspond to one or more descriptions of clinical trials sharing the same vaccine strategy. The output of this field detection step is a temporal series of clustering  $\mathcal{C}^* = \{C^T | T \in \mathcal{T}^*\}$  with  $C^T = \{C_j | j \in \mathcal{J}^T\}$  and  $C_j = \{r_i | r_i \in \mathcal{L}, i \in \mathcal{I}_j \subset \mathcal{I}\}$  computed over all the periods. It describes all the research directions explored from February 2020 to October 2021.

4. **Inter-temporal matching.** A temporal matching algorithm is then applied to identify meaningful kinship connections between fields from one period of time to another, *i.e.* fields that belong to the same research stream. We finally highlight the different research streams  $B_k$  over time and called them *branches of knowledge* (Figure 1.4).

The phylomemy reconstruction process makes it possible to draw the knowledge lineages at different resolutions through the tuning of a *level of observation* (1). The complexity of the resulting semantic landscape can range from a wide ‘continent’ to an ‘archipelago’ of specialized branches of knowledge.

**Visualizing phylomemies.** The structures highlighted by a phylomemy reconstruction process synthesize the complexity of the knowledge produced by a research community. In order to make this newly reconstructed knowledge actionable and explorable, a phylomemy can be visualized as a temporal network with time going by from top to bottom (8). Fields are represented by full circles and solid dark lines translate their kinship connections. *Emerging terms*<sup>§</sup> are displayed over the whole structure according to the combined coordinates of their period and fields of appearance. Term’s size depends on their frequencies in the original corpus of trials. Branches are sorted from left to right so that closely related ones lie side by side. Interactive features can be used to reveal the entire fields’ content, follow the dissemination of a given term throughout the phylomemy or simplify the *scale of description* of a selected branch.

<sup>‡</sup>The full list of roots is available at <https://doi.org/10.7910/DVN/JTR17A>.

<sup>§</sup>Terms appearing for the first time in the phylomemy

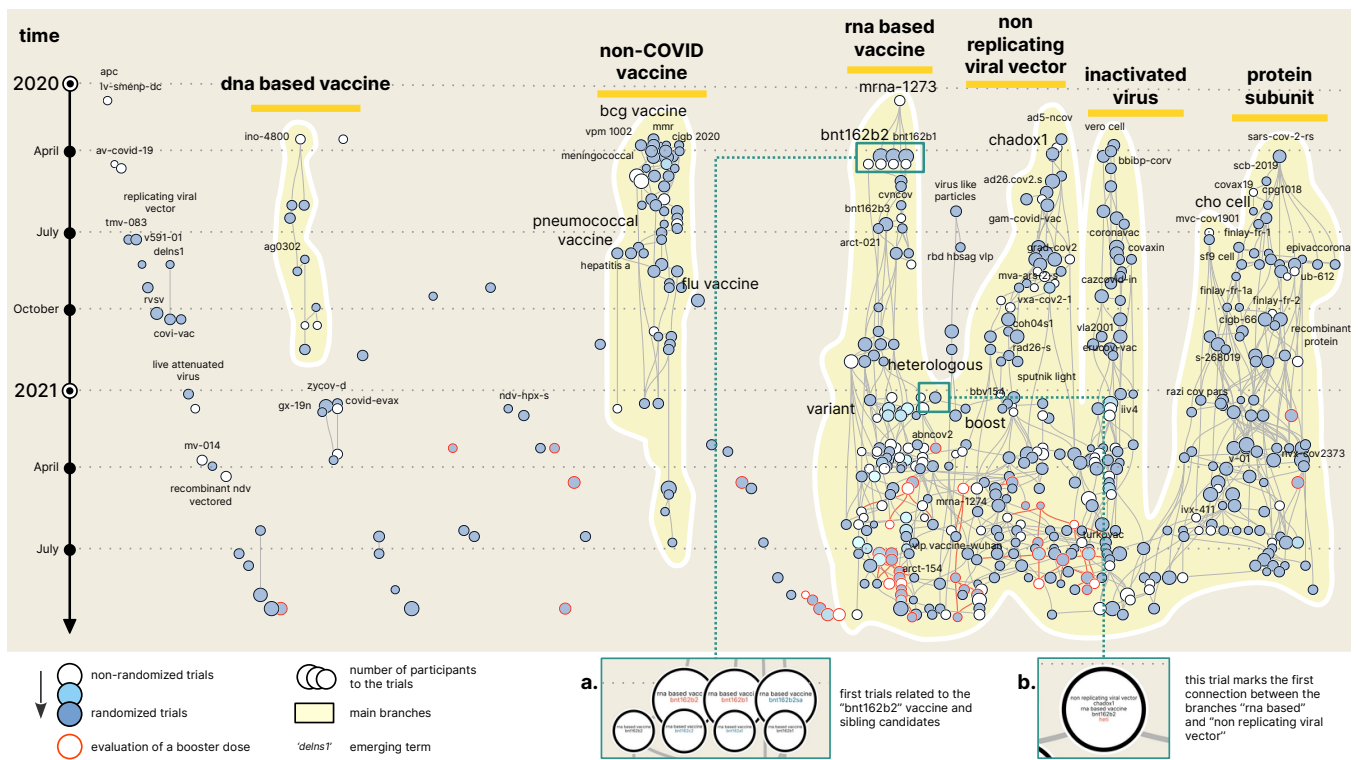
## Description of the resulting phylomemy

For our case study, we have used the corpus  $\mathcal{D}_{vt}$  of 1,794 COVID-19 vaccines clinical trials (see [Materials](#)) in order to reconstruct the weekly evolution of the research on COVID-19 vaccines between February 2020 and October 2021. Key expressions are extracted from the original descriptions of the tested vaccines, grouped into roots and then into fields. The reconstructed fields thus embody a set of trials at a given period of time. We here choose a level of observation  $\lambda = 0.5$  to shape quite precise branches. The resulting phylomemy (Figure 2) contains 175 roots and 550 fields distributed among 55 branches. The largest ones ‘dna based vaccine’, ‘non-COVID vaccines’, ‘rna based vaccine’, ‘non-replicating viral vector’, ‘inactivated virus’ and ‘protein sub-unit’ are highlighted by yellow shades. Shades of blue indicate the proportion of randomized clinical trials among the total number of trials on which the corresponding field has been reconstructed. The visualization of a phylomemy can also offer its user to interactively highlight some key information, as for example the research paths addressing vaccine boost issues, highlighted in red at the bottom of this figure.

## Following the worldwide tracks of COVID-19 vaccines

**General observations.** After having explored and analyzed Figure 2 alongside epidemiologists, we noticed that the reconstructed phylomemy clearly retrieves five major COVID-19 vaccine platforms in the form of complete branches. These platforms include the classical vaccine platforms *i.e.*, ‘non-replicating viral vector’, ‘inactivated virus’ and ‘protein sub-unit’ as well as the next-generation vaccine platform *i.e.*, ‘dna based vaccines’ and ‘rna based vaccines’. The visualization shows the continuous development of each branch and the way some of them started to interact and eventually blended while others stopped. Interestingly, trials of ‘rna based vaccines’ were registered very early in the course of the pandemic (February 2020) with trials evaluating the vaccine developed by Moderna TX (mRNA-1273) followed by the vaccine developed by Pfizer/BioNTech (BNT162b2) and sibling ones like BNT162b1 or BNT162b2sa that were not much longer tested (see Figure 2.a). The number of trials increased rapidly and interactions with other widely explored techniques were observed shortly afterwards: notably with the ‘non-replicating viral vector’ family (ChAdOx1 – AstraZeneca – see Figure 2.b). The latest interaction involved the ‘protein subunit’ branch in July 2021. In contrast, ‘dna based vaccines’, with a first trial registered in April 2020, had a very limited number of trials planned and the whole branch stopped rapidly in 2020. Similarly, other platforms of ‘replicating viral vector vaccine’, ‘virus-like particle vaccine’ and ‘live attenuated virus vaccine’ showed a very limited development.

**Repurposing non-COVID vaccines.** As the development and approval of COVID-19 vaccines was expected to take time, researchers also explored repurposing non-COVID vaccines. Considering the lower severity of the disease in children and young adults, some researchers hypothesized the possible heterologous protective effect of these vaccines. Some evidence shows that live-attenuated vaccines such as Bacille Calmette–Guerin (BCG), Measles, Mumps, Rubella (MMR) can induce protective innate immunity, which could be central in controlling SARS-CoV-2 (9). While this hypothesis was appealing, it



**Fig. 2.** Phylogenomy of 1,794 COVID-19 vaccines trials recorded between February 2020 and October 2021 in the COVID-NMA database. Online and interactive version available at [maps.gargantext.org/publications](https://maps.gargantext.org/publications)

165 did not seem to expand into a wider research domain. The  
 166 the branch of ‘non-COVID vaccines’ appears and expands at the  
 167 beginning of the pandemic but progressively decreases towards  
 168 the end of 2020 as other more promising vaccines arose. Nev-  
 169 ertheless, some researchers highlighted the need to adequately  
 170 assess the use of non-COVID live-attenuated vaccines as they  
 171 could potentially boost response in high-risk populations, be  
 172 used in addition to COVID-vaccines to increase effectiveness  
 173 and durability of their effect, or be used to protect people  
 174 exposed to COVID-19 patients (9).

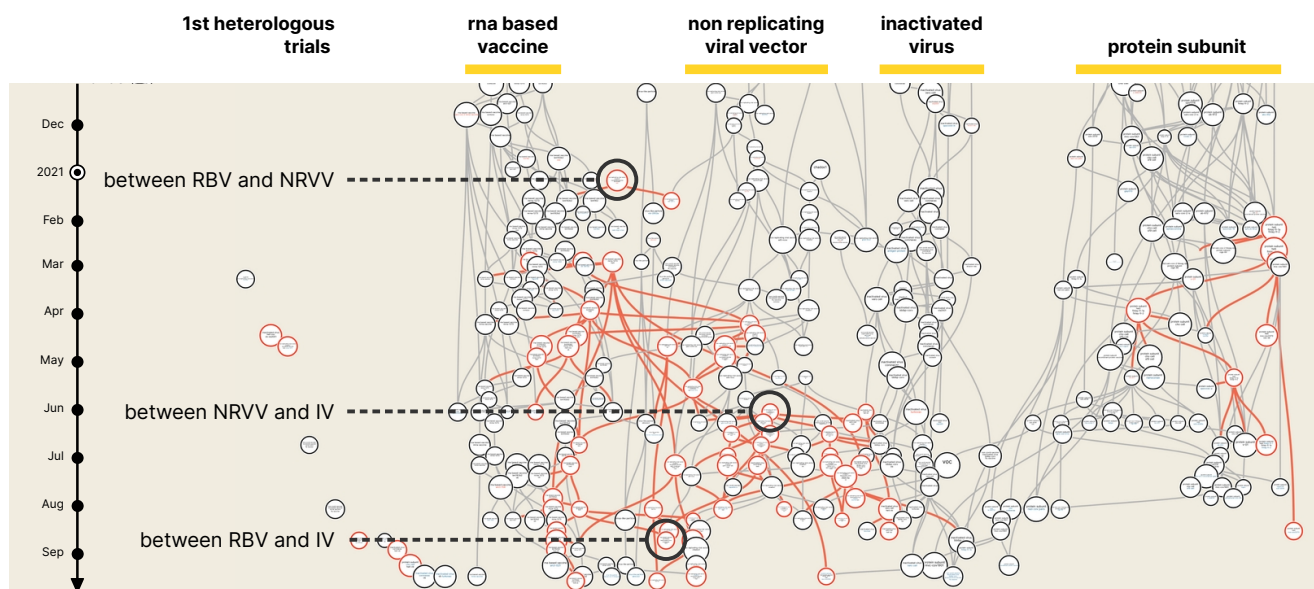
175 **Heterologous vaccination.** The branches interactions reflect  
 176 the exploration of a new approach to vaccine implementation  
 177 moving from homologous prime vaccination (i.e., injections of  
 178 two doses of the same vaccine) to heterologous prime vacci-  
 179 nation (i.e., injection of the first dose of a given vaccine and  
 180 the second dose of another vaccine). This is clearly shown in  
 181 Figure 3 with the assessment of the heterologous prime vaccina-  
 182 tion of ‘rna based vaccine’ (BNT162b2-Pfizer/BioNTech) and  
 183 ‘non-replicating viral vector’ (ChAdOx1-AstraZeneca) in early  
 184 2021. This new approach was motivated by concerns about  
 185 waning vaccine immunity, but also by practical considerations.  
 186 Following concerns about the safety of the AstraZeneca ChA-  
 187 dOx1 vaccine, the EMA recommended giving a second dose  
 188 Pfizer BNT162b2 vaccine to patients under the age of 55 years  
 189 old who received one dose of ChAdOx1-S-nCoV-19. Further-  
 190 more, decision makers needed flexibility to overcome the issue  
 191 of vaccine availabilities during the vaccine rollout. This new  
 192 approach proved to be relevant and other associations were  
 193 evaluated: ‘non replicating viral vector’ and ‘inactivated virus’  
 194 in June 2021 and later ‘rna based vaccine’ and ‘inactivated  
 195 virus’ in September 2021.

**Boosters.** Phylogenemies are essential in identifying shifts in  
 196 research questions. While evidence of the beneficial effect of  
 197 vaccines is mounting, research questions are moving toward  
 198 exploring the effect of booster to overcome the waning of vac-  
 199 cine efficacy over time. Early in 2021, new trials assessing the  
 200 impact of administrating a third dose (see Figure 2, red outline  
 201 at the bottom) have been registered particularly for ‘rna based  
 202 vaccines’ and ‘non-replicating viral vector’ (10). An impor-  
 203 tant part of the research on boosters’ effects is considering  
 204 heterologous boosters.  
 205

**Filters and upcoming research questions.** By using additional  
 206 data from the trials registries, we can filter the current phy-  
 207 lomemy and thus push faceted observations to the fore or  
 208 identify upcoming research questions.  
 209

Phylogenemies also provides important information on re-  
 210 search planning and reporting. As shown in Figure 2, most  
 211 trials registered are randomized controlled trials. Early in the  
 212 pandemic, non-randomized trials were primarily early phase  
 213 trials while those registered in 2021 include both early phase  
 214 trials exploring new vaccines and phase 4 trials assessing vac-  
 215 cines safety.  
 216

We can also explore the visualization to better under-  
 217 stand how different countries participated in the overall re-  
 218 search effort over time. For example, when filtering on the  
 219 country (see [maps.gargantext.org/phylo/vaccines/countries](https://maps.gargantext.org/phylo/vaccines/countries)),  
 220 we see that trials conducted in the USA explored all vac-  
 221 cine platforms and that first registered trials frequently in-  
 222 volved a center in the USA, confirming their leading role in  
 223 clinical research (e.g., ‘dna based vaccine’, ‘rna based vac-  
 224 cine’, ‘protein subunit’). Other important trials character-  
 225 istics such as funding sources can also be highlighted (see  
 226



**Fig. 3.** A focus of Figure 2. In red are highlighted all the trials evaluating heterologous primary vaccination and heterologous booster. We circle the first heterologous trials involving different platforms.

227 [maps.gargantext.org/phylo/vaccines/fundings](https://maps.gargantext.org/phylo/vaccines/fundings)).

228 Finally, we address the question of the publication of trial  
 229 results (i.e., preprint or peer-reviewed articles). As shown in  
 230 Figure 4, we currently have access to the results of a very  
 231 limited number of planned trials. While most of the COVID  
 232 vaccine trials registered in early 2020 are published, most of  
 233 the non-COVID vaccine trials are still unpublished. Under-  
 234 standing whether these trials were actually conducted with  
 235 unpublished results or were unable to recruit is an impor-  
 236 tant issue.

### 237 Perspectives and insights for COVID-19 research

238 Global coordination between research teams is a key for accel-  
 239 erating innovation in Science, especially during crisis situations  
 240 such as the COVID-19 pandemics. Reducing redundancies  
 241 and providing heuristics to find new search paths as they  
 242 arise can save time and lives (3). We claim that phylomemy  
 243 reconstruction could be instrumental to guide trialists, fund-  
 244 ers and decision makers in biomedical research. In times of  
 245 crises, it would enable them to better adapt to the evolution  
 246 of the situation by following emerging research questions and  
 247 identify less promising domains. It could also facilitate the  
 248 identification of research gaps, research questions that may  
 249 have been abandoned prematurely and redundancy in research.  
 250 Our phylomemies could also be enriched with other data :

- 251 • data already recorded in the trials registries such as out-  
 252 comes or participants characteristics which would allow  
 253 exploring research conducted on vulnerable populations  
 254 (children, pregnant women, immunocompromised patients,  
 255 elderly etc.), trial results posted on the registries when  
 256 possible;
- 257 • data that are not part of the registries but which should  
 258 be added in pandemic times like the number of patients  
 259 actually included in the trials;

- data that exists outside of the registries (publications,  
 trials results, etc.) but for which a difficult work of data  
 pruning and integration is required.

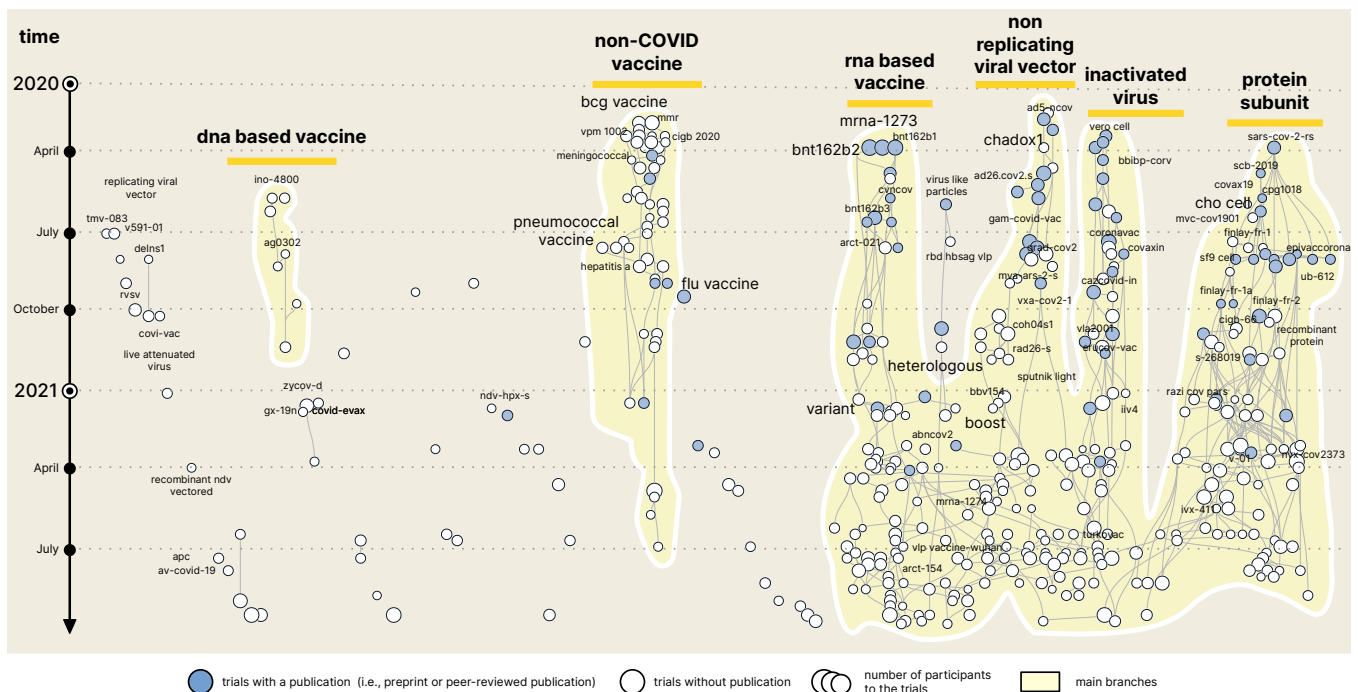
The addition of such information could be fulfilled through  
 the collaborative and cumulative features of the *Gargantext*  
 platform<sup>1</sup>: the software used for computing the phylomemy  
 reconstruction process. This generalization would increase  
 the benefits of this approach tenfold. In a world where ex-  
 perts are increasingly specialized, it could draw attention to  
 alternative solutions developed in other branches of science  
 or to problems already encountered in research direction to  
 be explored. It could also lead to new conceptual operations  
 to be performed on a knowledge database, such as "give me  
 all the branches of knowledge that are merging" or "suggest  
 a promising combination of compounds to test". This could  
 both accelerate research by making tangible the latent struc-  
 ture of innovation, and promote collaborations between teams  
 that would not otherwise be interested in each other's work.  
 Phylomemy reconstructions may thus become collective and  
 reflective tools to foster the worldwide collective coordination  
 between researchers. This revolution in clinical trial processing  
 is within reach. Nevertheless, it would imply having access to  
 high quality data on research planning and protocol.

Our case study focuses on a single disease, but this approach  
 is fully generic and we call for a worldwide observatory for  
 monitoring the dynamics of clinical trials. As it scales up, our  
 approach could be implemented for any disease or research  
 field.

### Materials

Our data set has been collected and curated by the combined effort of  
 epidemiologists, data integration and complex systems researchers.

<sup>1</sup>Gargantext is a free software. See <https://iscpi.fr/projects/gargantext>



**Fig. 4.** Phylomemy of the randomized only COVID-19 vaccines trials. In blue, we highlight all the trials with an associated publication (i.e., preprint or peer-reviewed articles).

291 **The COVID-NMA database.** The COVID-NMA project is an international  
 292 initiative aimed at providing a living mapping and a living  
 293 systematic review of all trials assessing treatments and preventive  
 294 interventions for COVID-19 (2, 3). The development of the  
 295 COVID-NMA database relies on a full methodology designed to  
 296 generate and make available a complete, comprehensive, integrated,  
 297 non-redundant and carefully annotated data sets on clinical trials.  
 298 We automatically extract data from clinical registries on a weekly  
 299 basis and provide assistance to epidemiologists on the curation and  
 300 annotation process. Raw data is extracted from the [EU clinical  
 301 trials register](#), from the [ClinicalTrial registry](#) managed by the U.S.  
 302 National Library of Medicine, from the [IRCT registry](#) and from the  
 303 [WHO International Clinical Trials Registry Platform \(ICTRP\)](#)  
 304 – an international registry that assembles information on clinical  
 305 trials registered in 17 primary registries to identify new trial  
 306 assessing COVID-19 vaccine and update of previously registered trial  
 307 records. Data are extracted from registries, annotated by epidemi-  
 308 ologists, then stored and made available through the COVID-NMA  
 309 database<sup>||</sup>.

310 **Pre-processing the database.** We have pre-processed\*\* the COVID-  
 311 NMA database before using it for the phylomemy reconstruction  
 312 to filter the 1,794 descriptions related to vaccines trials. The trials  
 313 records have been first aggregated by publication week. Then, we  
 314 have merged the sections ‘*pharmacological treatment*’, ‘*treatment  
 315 type*’ and ‘*treatment name*’ together to shape the trial descriptions.  
 316 These descriptions have also been enriched with extra-information  
 317 such as trial phases, funding, involved countries or associated pub-  
 318 lications. The resulting corpus  $\mathcal{D}_{vt}$  has latter been collectively  
 319 and collaboratively curated by epidemiologists thanks to the free  
 320 software *Gargantext* (5). There, these experts have extracted and  
 321 validated a core vocabulary as a list of 175 root terms.

322 **Data Availability.** The original COVID-NMA database  
 323 can be downloaded at [covid-nma.com](#). The recon-  
 324 structed phylomemy is available for live explorations at

[maps.gargantext.org/phylo/vaccines/publications](https://maps.gargantext.org/phylo/vaccines/publications) and download-  
 able at <https://doi.org/10.7910/DVN/JTR17A>.

1. D Chavalarias, Q Lobbé, A Delanoë, Draw me science – multi-level and multi-scale recon- 327  
 struction of knowledge dynamics with phylomemes. *Scientometrics* (2021). 328
2. I Boutron, et al., The COVID-NMA Project: Building an Evidence Ecosystem for the COVID- 329  
 19 Pandemic. *Ann Intern Med* **173**, 1015–1017 (2020). 330
3. VT Nguyen, et al., Research response to coronavirus disease 2019 needed better coordina- 331  
 tion and collaboration: a living mapping of registered trials. *J Clin Epidemiol* **130**, 107–116 332  
 (2021). 333
4. D Chavalarias, JP Cointet, Phylomemetic patterns in science evolution—the rise and fall of 334  
 scientific fields. *PLoS one* **8**, e54847 (2013) 00000 bibtex: chavalariasPhylomemetic2013. 335
5. A Delanoë, D Chavalarias, Mining the digital society - Gargantext, a microscope for collabo- 336  
 rative analysis and exploration of textual corpora. (forthcoming 2021). 337
6. G Dias, R Mukelov, G Cleuziou, Mapping general-specific noun relationships to wordnet 338  
 hypernym/hyponym relations in *International Conference on Knowledge Engineering and* 339  
*Knowledge Management*. (Springer), pp. 198–212 (2008). 340
7. T Uno, M Kiyomi, H Arimura, , et al., Lcm ver. 2: Efficient mining algorithms for fre- 341  
 quent/closed/maximal itemsets in *Fimi*. Vol. 126, (2004). 342
8. Q Lobbé, A Delanoë, D Chavalarias, Exploring, browsing and interacting with multi-level and 343  
 multi-scale dynamics of knowledge. *Inf. Vis.*, 14738716211044829 (2021). 344
9. K Chumakov, et al., Old vaccines for new infections: Exploiting innate immunity to control 345  
 covid-19 and prevent future pandemics. *Proc. Natl. Acad. Sci.* **118** (2021). 346
10. PR Krause, et al., Considerations in boosting COVID-19 vaccine immune responses. *Lancet* 347  
**398**, 1377–1380 (2021). 348

<sup>||</sup>We here note that international trials registries can be post-updated by research teams, e.g. for post-adding a related publication. Future versions of the phylomemies presented in this paper might thus be slightly different from the current ones. A promising way to get around this issue would be to archive every modifications of the original registries and then choose the version we want to integrate in the phylomemies.

\*\*The pre-processing script can be downloaded at <https://doi.org/10.7910/DVN/JTR17A>