

822 **A Further analysis of Fedivertex**

823 This appendix presents several additional plots complementary to the figures presented in the main
 824 text. In particular, we propose to further analyze the degree distribution of our graphs, and the
 825 language distribution within different Fediverse social media.

826 **Degree distribution** Fig. 5 and 6 details the degree distribution for each Fediverse software.

827 For federation graphs, we identify two distinct patterns. First, in BookWyrM, Lemmy, and Friendica,
 828 the degree distribution exhibits a peak at high degrees, and the proportion of nodes decreases with the
 829 degree. This behavior is particularly pronounced in Lemmy. In contrast, Mastodon, Misskey, and
 830 Pleroma exhibit a peak at low degrees (close to 0), followed by a relatively flat region for intermediate
 831 degrees, and a sharp drop for high degrees. This similarity can be attributed to the fact that these
 832 platforms are all centered on the same activity: micro-blogging.

833 In the active user graphs, the degree distribution consistently follows a power-law-like pattern, with
 834 the number of nodes decreasing as the degree increases. The intra-instance graph of Lemmy exhibits
 835 a similar trend, though the pattern is less pronounced.

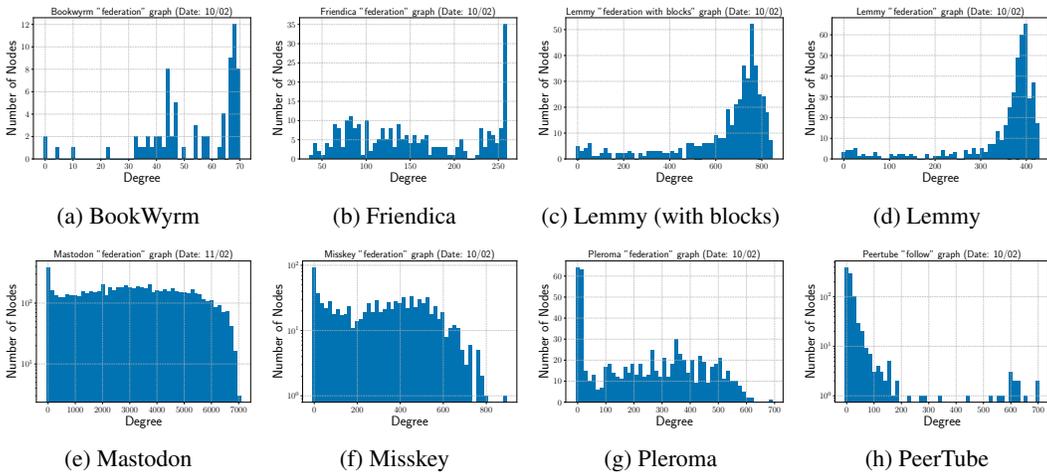


Figure 5: Degree-histogram distributions for federation graphs. Version with and without block for Lemmy, and Follow graph for Peertube

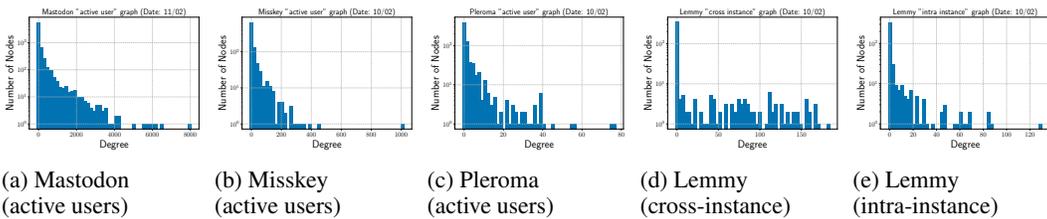


Figure 6: Degree-histogram distributions for the three active-user graphs and Lemmy’s cross- vs. intra-instance graphs.

836 Fig. 7 reproduces fig. 2 from the main text, but with a normalized degree: the degree is divided by the
 837 number of nodes. This normalization incorporates graph density: the closer a curve is to the top-right
 838 corner, the denser the graph. From fig. 7, Fedivertex graphs appear denser than the SNAP social
 839 graphs. While the two groups overlap in fig. 2, they are more distinctly separated in the normalized
 840 plot. This further illustrates major differences between Fedivertex and widely-used datasets.

841 **Language distribution** Fig. 8 presents histograms of the language distribution in Lemmy, Peertube,
 842 Mastodon, and Misskey. This figure shows diverse language distributions among Fediverse social
 843 media. First, Mastodon and Lemmy have a large majority of English-speaking instances. Second,

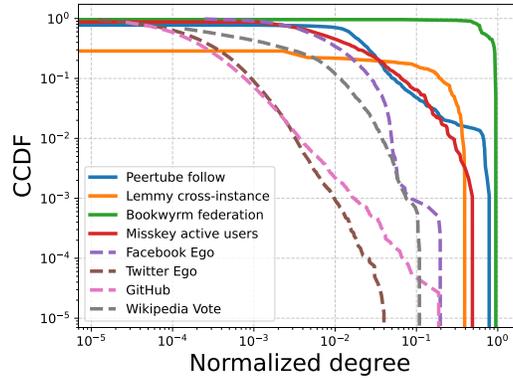


Figure 7: Complementary Cumulative Distribution Function (CCDF) of the degree for several Fediverse graphs and other widely-used graphs [34] based on social networks after normalization by the total number of nodes

844 Peertube is dominated by European languages. We observe more diversity than in Mastodon because
 845 French and German have a number of nodes similar to English. Third, Misskey shows a unique
 846 language distribution within the Fediverse as the Asian languages (especially Japanese) are the
 847 most spoken languages. This phenomenon is explained by the fact that Misskey (contrary to other
 848 Fediverse software) has been developed by a Japanese developer.

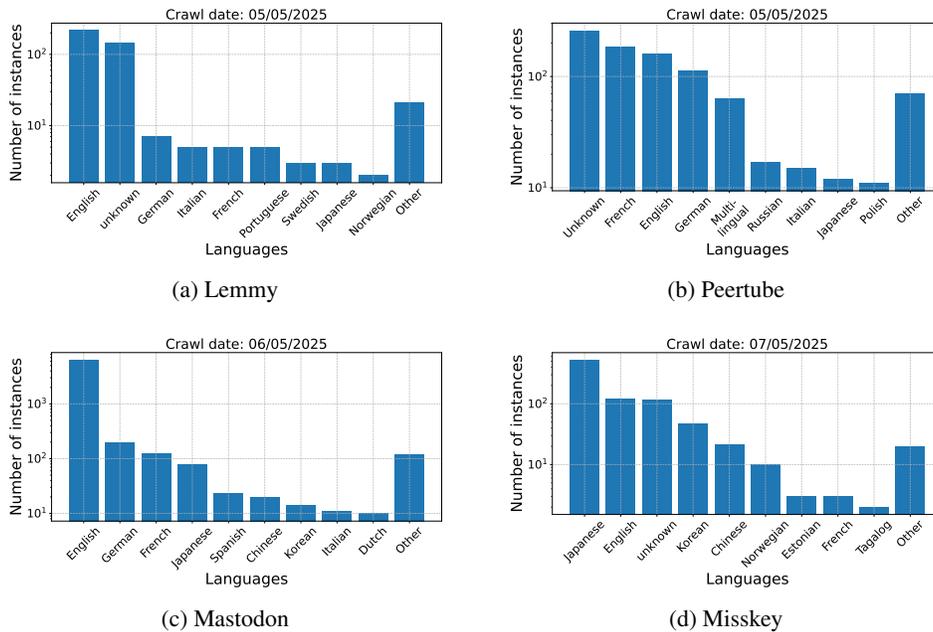


Figure 8: Histograms of the language distribution among four different Fediverse software, logarithmic scale

849 **For a thorough exploratory data analysis** we refer to our 5 Kaggle notebooks that cover these
 850 aspects, but also other dimensions such as weight or user distribution:
 851 <https://www.kaggle.com/datasets/marcdamie/fediverse-graph-dataset/code>.

852 B Python API examples

853 Our Python package `Fedivertex` produces a straightforward interface to interact with our dataset.
 854 This package seamlessly downloads and extracts graphs from our CSV files using the `NetworkX`

```

from fedivertex import GraphLoader

loader = GraphLoader()
G = loader.get_graph(software = "misskey", graph_type = "active_user")

nb_nodes = G.number_of_nodes()
nb_edges = G.number_of_edges()
degree = G.degree()

```

Listing 3: Code example to extract statistics about a Fedivertex graph

```

from fedivertex import GraphLoader
loader = GraphLoader()

# Extract the latest graph
G_latest = loader.get_graph(software = "misskey", graph_type = "active_user")

# Extract the i-th graph in the dataset (sorted chronologically)
i = 3
G_i = loader.get_graph(software = "misskey", graph_type = "active_user", index=i)

# Extract the graph of a specific data
G_2403 = loader.get_graph(software = "misskey", graph_type = "active_user", date="20250324")

```

Listing 4: Code example of the graph selection in Fedivertex

855 format, a popular Python package to manipulate graphs. This appendix demonstrates several simple
856 operations that can be performed using our package.

857 The installation is straightforward as our package is available on the Python Package Index (PyPI):
858 `pip3 install fedivertex`. The dataset is downloaded automatically the first time a graph
859 is loaded. The dataset is stored in cache thanks to the CroissantML package in the folder
860 `.cache/croissant`.

861 Listing 3 shows how to extract basic statistics from the latest Misskey active user graph. Our method
862 `get_graph` outputs a NetworkX object, so we can apply all functions and methods from NetworkX.
863 This example simply calls `number_of_nodes`, `number_of_edges`, and `degree`, but one can use
864 the entire NetworkX API. We refer to NetworkX documentation for further information about the
865 capabilities of their API: <https://networkx.org/documentation/stable/index.html>

866 In Fedivertex, a graph is defined using three elements: the software (e.g., Misskey), the graph type
867 (e.g., active user), and the date. For date selection, the package provides three different options,
868 illustrated in listing 4. First, if no date is provided to `get_graph`, the latest graph is automatically
869 selected. Second, if the user provides an (integer) index i , the i -th graph is selected. This index works
870 like a Python list index, so -1 is accepted and the indexing starts at 0. Third, the user can provide a
871 date in a string format YYYYMMDD (e.g., 20250324).

872 As our dataset contains a large range of possible software, graph types, and dates, we also
873 provide utility functions to simplify the interactions. Listing 5 presents these functions. First,
874 `list_all_software` lists all the available software. Second, `list_graph_types` lists the avail-
875 able graph types for a given software. Third, `list_available_dates` lists the available dates for a
876 given software and graph type.

```

from fedivertex import GraphLoader
loader = GraphLoader()

print(loader.list_all_software())
# ["bookwyrm", "friendica", "lemmy", "mastodon", "misskey", "peertube", "pleroma"]

print(loader.list_graph_types("peertube"))
# ["follow"]

print(loader.list_available_dates("peertube", "follow"))
# ["20250203", "20250210", "20250217", "20250224", "20250303", "20250311", "20250317",
# "20250324", "20250331", "20250407", "20250414", "20250421", "20250428", "20250505"]

```

Listing 5: Code example of the utility functions in Fedivertex

877 As seen in listing 1 in main text, it is possible to use graph extraction with the option
 878 “only_largest_component = True.” When set to true, the method only returns the largest compo-
 879 nent of the graph. As many graph algorithms require connected graphs, this option eases benchmarks
 880 of such algorithms. Additionally, listing 2 shows how to extract the language information to use it as
 881 a ground-truth label (e.g., in community detection like in table 4)

882 For further advanced examples, we refer to our Kaggle notebooks that analyze the dataset using the
 883 Python package:
 884 <https://www.kaggle.com/datasets/marcdamie/fediverse-graph-dataset/code>.

885 C Experiments on defederation

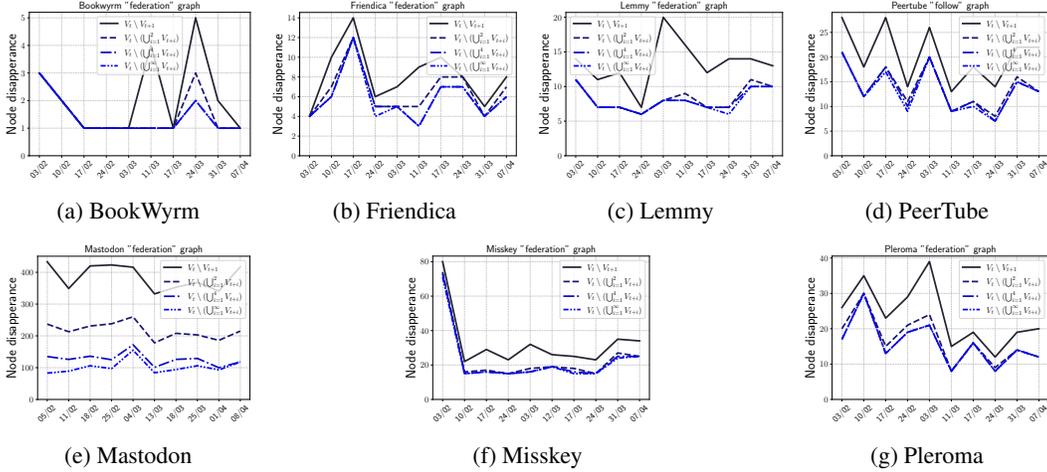


Figure 9: Node deletion at different time horizons for federation graphs, and the follow graph (for PeerTube).

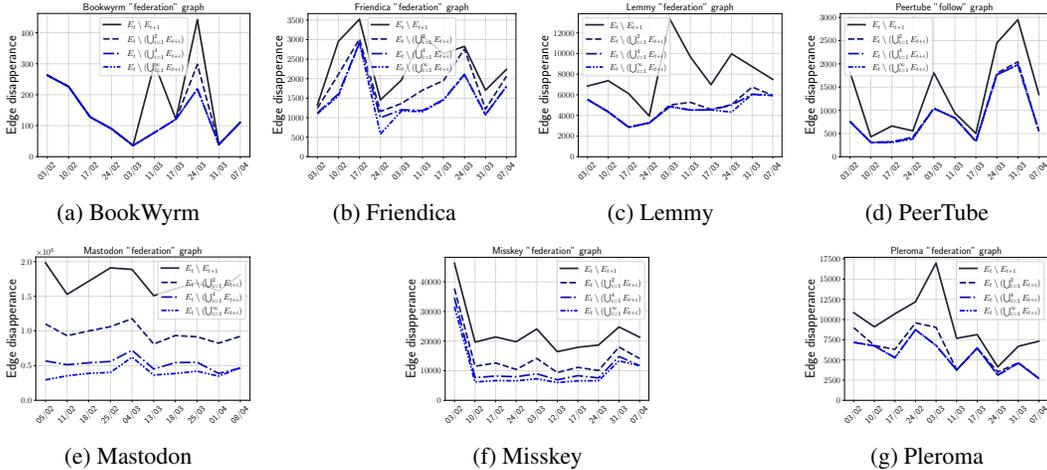


Figure 10: Edge deletion at different time horizons for federation graphs, and the follow graph (for PeerTube).

886 In section 5.2, we focused on edge deletion prediction by using the same method than for edge
 887 prediction, with limited results. A natural question is thus: is the poor performance due to noise in the
 888 deletion process? To answer this question, we must assess whether edge deletion is a random process
 889 (e.g., caused by network instability) or whether it persists over time. In particular, fig. 4 does not
 890 exhibit clear trends and tends to suggest that deletions could be accidental. We test this hypothesis
 891 by comparing the edges and nodes deleted between two consecutive weeks (as reported in fig. 4) to

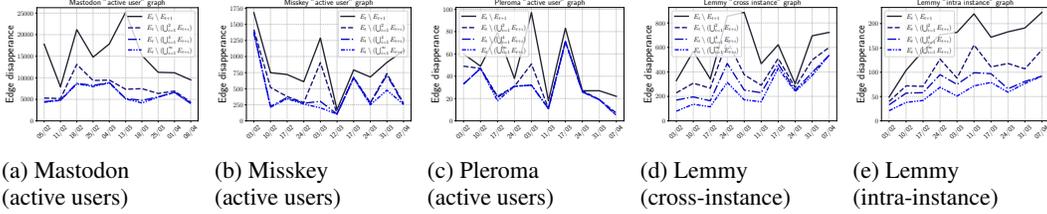


Figure 11: Edge deletion at different time horizons for the three active-user graphs and Lemmy’s cross- vs. intra-instance graphs.

Graph	Betweenness	Eigenvector centrality	Pagerank	Random
Misskey AU	6	16	16	7 ± 2.4
Misskey fed.	7	12	11	7 ± 2.3
Friendica fed.	11	13	13	7 ± 2.4
Pleroma fed	13	3	12	8 ± 2.5

Table 5: Comparison of Betweenness, eigenvector centrality and Pagerank scores on different graphs by reporting the number of correct predictions in Top-50 scores (the higher the better)

892 those that remain deleted over longer time horizons: two weeks, four weeks, and until the end of the
 893 measurement period.

894 We report the results for node deletion in fig. 9, and for edge deletion in fig. 10 and fig. 11, depending
 895 on the type of graph. The results show that most variations are permanent, at least within the time
 896 scale of our study, as soon as they persist beyond two weeks.

897 The majority of deletions appear to be permanent, except in BookWorm, which can be explained
 898 by the small size of the network, and in Mastodon, where node and edge deletion durations follow
 899 a continuum. The only other case showing a progressive increase in the proportion of affected
 900 edges with respect to duration is Lemmy cross- and intra-instance, which can be explained by the
 901 construction of these graphs: they depend on recent message exchanges between instances and may
 902 naturally include servers whose users interact only occasionally.

903 In these cases, edge deletion does not appear to be an artifact of the crawling process, but rather
 904 reflects actual communication dynamics in the graph.

905 We also present a baseline for node deletion. Similarly to the case of edges, we report the number of
 906 actual deletions among the top-50 nodes with the highest scores. We propose to predict deletion for
 907 nodes that are the least central in the graph, as this may indicate lower activity. Other metrics, and
 908 incorporating additional information about the nodes, could certainly improve these predictions. We
 909 report the results in table 5.

910 D Additional remarks about the crawler

911 We discuss in this section the robustness of the crawling, and legal procedure prior the crawler
 912 deployment. We refer to section 3.5 for the ethical concerns related to the crawling procedure.

913 **Robustness** Fig. 4c represents the variations of the number of nodes in different Fediverse software.
 914 For all these social networks, the variations looks erratic. Such patterns could raise concerns regarding
 915 the robustness of the crawling procedure. However, these erratic variations are a natural phenomenon
 916 of the Fediverse and are not an artifact created by the crawler.

917 To put these variations into perspective, fig. 12 represents the overall evolution instead of the node
 918 variation between two snapshots. We observe that the number of nodes is relatively stable. The
 919 patterns from fig. 4c are barely distinguishable on fig. 12 because they concern a minority of nodes.

920 Even though these patterns involve only a minority of nodes, they are worth examining. Most
 921 instances on the Fediverse are maintained by volunteers who cover the operational and maintenance
 922 costs themselves. Many are operated by a single individual who may lack the time, financial resources,

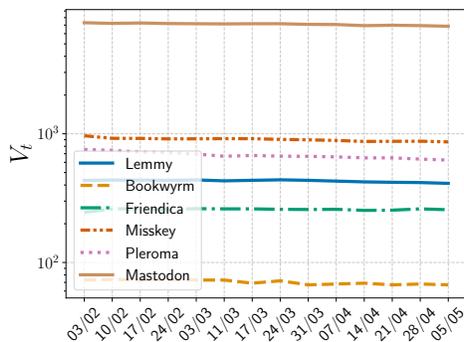


Figure 12: Evolution of the number of nodes over time for the six software of Fediverse

923 or motivation to ensure long-term stability. While the biggest instances benefit from a high stability
 924 due to their teams of experienced engineers, smaller nodes can quickly appear, disappear, and
 925 eventually reappear.

926 A natural question that follows from this observation is whether our crawler adequately captures
 927 the dynamics of unstable instances. Our approach relies on a curated list maintained by Fediverse
 928 Observer, which continuously crawls the Fediverse to discover new instances and assess the availabil-
 929 ity of existing ones. Notably, instances that are temporarily unavailable remain listed, meaning our
 930 crawler does not ignore unstable nodes. Thanks to the frequent updates and broad coverage provided
 931 by Fediverse Observer, the resulting instance list is both extensive and up-to-date. We thus argue that
 932 our crawler offers a realistic and consistent view of the Fediverse.

933 Our source code is publicly available: <https://github.com/MarcTOK/Franck>. The crawler
 934 implementation has remained consistent since the first crawl, with only minor changes introduced to
 935 improve logging completeness.

936 **Legal compliance** Before crawler deployment, we have been involved in active discussions with
 937 the legal departments of our institutions to cover any potential legal issue. We worked in particular
 938 on the compliance with data privacy regulations, especially GDPR. As stressed in section 3.5, our
 939 graphs only contain aggregated information (i.e., server-level information), and we only query public
 940 APIs, and thus are less privacy sensitive than existing works (e.g., the Webis Mastodon corpus [55]
 941 that compiled 700 million Mastodon posts). The crawling logs are kept only for the time necessary to
 942 confirm the crawl success and deleted afterwards.

943 Our work is among the first research projects on the Fediverse, and the guidelines for research in this
 944 field still need to be refined. Existing guidelines tend to focus on centralized mainstream social media,
 945 and are thus not well adapted to scenarios where instances are run by distinct legal entities, possibly
 946 subject to different regulations and norms. The discussions with the legal departments allowed
 947 us to clarify the legal requirements in this context. We shared our experience with the Network
 948 of Alternative Social Media Researchers, a transdisciplinary international collective of academics
 949 studying the cultures, technologies, and practices of non-corporate social media⁵, to help formalize a
 950 dedicated ethical framework for research on alternative social media.

951 In addition to the discussions with the legal department, we have also coordinated our efforts with
 952 the IT department to avoid any disturbance caused by our crawler. Notably, our crawler is formally
 953 declared to the security operations center of our university, is assigned specific IP addresses, and is
 954 under the surveillance of the IT department. These measures are complementary to the preventive
 955 strategy of a *slow* crawl, aiming to minimize disturbances for the instances, as detailed in section 3.5.
 956 After four months of experiments, we have received no complaints, neither from our university nor
 957 from Fediverse instance administrators.

⁵Testimony at <https://www.socialmediaalternatives.org/2025/05/07/asm-research-ethics.html>