

Torchlight: Diffusion-based Network Trace Generation from the DARPA Searchlight Dataset

RAY ZHAO, University of Southern California,
ALEFIYA HUSSAIN, USC/ISI,

At present, there is a severe lack of both comprehensive and realistic labeled datasets for machine learning applications in the networking domain. Predominantly, prior generative work has focused on lower-dimensional representations that rely on aggregating flow characteristics and lack the fine-grain of raw network traces. This results in suboptimal performance in machine learning contexts and limited applications outside of those contexts. This has induced a push for new generative techniques to provide synthetic data for usage both on its own and layered in with real data as augmentation. In this paper, we present Torchlight, a diffusion-based generation framework built atop and extending techniques first introduced in *NetDiffusion* [?] using DARPA Searchlight [?] data to generate synthetic network traces for video streaming applications. We demonstrate the efficacy of Torchlight in generating synthetic network traces that reasonably resemble real-world data and perform notably well in classification tasks.

1 INTRODUCTION

There's a large need and general scarcity of labeled network datasets, as high-quality data is often expensive and fraught with privacy concerns, and the datasets that do exist are often left rarely update over time. Synthetic data generation techniques aim to solve this problem, but prior GAN based methodologies have been limited in their ability to generate raw, high-dimensional traces and often rely on aggregating flow characteristics, which results in data that lacks statistical fidelity to real data and performs much more poorly in machine learning tasks. The highly specific formatting of these types of condensed traces also renders them generally incompatible with traditional tools like Wireshark and tcpdump.

In recent years, the advent of diffusion models have solved many of the problems that plague older GAN frameworks, being able to capture much more complex patterns and relationships. They also usually offer much more stability in the training process, in which GANs are notably finicky. Additionally, their now near-hegemony in the image generation space has enabled a wide base of development support in the diffusion community and has also made available many "plug and play" means of constraining generation to appropriate enough scopes for specific generation tasks like the one at hand.

As it pertains to Torchlight, we build on the NetDiffusion framework pioneered by Jiang et al. [?] to generate synthetic network traces for video streaming applications, using the DARPA Searchlight dataset [?] as our base. Our contributions are as follows:

- (1) **Successful application of NetDiffusion techniques on the Searchlight Dataset:** Using the pre-processing and post-processing code provided by Jiang et al., we were able to successfully generate synthetic traces that performed capably on its own as training data for classification tasks as well as when used in conjunction with real data as augmentation.
- (2) **Streamlining the generation process:** While the original NetDiffusion codebase used externally-built WebUIs for LoRA fine-tuning as well as the Controlnet, these two major components were integrated into the codebase such that the entire process can be run from a single script.

2 MOTIVATION & RELATED WORK

Publicly available network data is crucial in this field for enabling continued research and development of new networking technologies and techniques. While useful for more traditional networking analysis tasks, machine learning applications have been driving much of the development in use-cases like anomaly detection and traffic classification and in particular have largely relied on these datasets.

2.1 Cost and Scarcity

The primary parameter given to the “acmart” document class is the *template style* which corresponds to the kind of publication or SIG publishing the work. This parameter is enclosed in square brackets and is a part of the `\documentclass` command:

```
\documentclass[STYLE]{acmart}
```

Journals use one of three template styles. All but three ACM journals use the `acmsmall` template style:

- `acmsmall`: The default journal template style.
- `acmlarge`: Used by JOCCH and TAP.
- `acmtog`: Used by TOG.

The majority of conference proceedings documentation will use the `acmconf` template style.

- `sigconf`: The default proceedings template style.
- `sigchi`: Used for SIGCHI conference articles.
- `sigplan`: Used for SIGPLAN conference articles.

2.2 Template Parameters

In addition to specifying the *template style* to be used in formatting your work, there are a number of *template parameters* which modify some part of the applied template style. A complete list of these parameters can be found in the *L^AT_EX User’s Guide*.

Frequently-used parameters, or combinations of parameters, include:

- `anonymous, review`: Suitable for a “double-anonymous” conference submission. Anonymizes the work and includes line numbers. Use with the `\print` command to print the submission’s unique ID on each page of the work.
- `authorversion`: Produces a version of the work suitable for posting by the author.
- `screen`: Produces colored hyperlinks.

This document uses the following string as the first command in the source file:

```
\documentclass[acmsmall]{acmart}
```

3 METHODS

Modifying the template — including but not limited to: adjusting margins, typeface sizes, line spacing, paragraph and list definitions, and the use of the `\vspace` command to manually adjust the vertical spacing between elements of your work — is not allowed.

4 EVALUATION

The “acmart” document class includes the “booktabs” package — <https://ctan.org/pkg/booktabs> — for preparing high-quality tables.

Table captions are placed *above* the table.

Because tables cannot be split across pages, the best placement for them is typically the top of the page nearest their initial cite. To ensure this proper “floating” placement of tables, use the environment **table** to enclose the table’s contents and the table caption. The contents of the

Table 1. Frequency of Special Characters

Non-English or Math	Frequency	Comments
Ø	1 in 1,000	For Swedish names
π	1 in 5	Common in math
\$	4 in 5	Used in business
Ψ_1^2	1 in 40,000	Unexplained usage

Table 2. Some Typical Commands

Command	A Number	Comments
<code>\author</code>	100	Author
<code>\table</code>	300	For tables
<code>\table*</code>	400	For wider tables

table itself must go in the **tabular** environment, to be aligned properly in rows and columns, with the desired horizontal and vertical rules. Again, detailed instructions on **tabular** material are found in the *TEX User’s Guide*.

Immediately following this sentence is the point at which Table ?? is included in the input file; compare the placement of the table here with the table in the printed output of this document.

To set a wider table, which takes up the whole width of the page’s live area, use the environment **table*** to enclose the table’s contents and the table caption. As with a single-column table, this wide table will “float” to a location deemed more desirable. Immediately following this sentence is the point at which Table ?? is included in the input file; again, it is instructive to compare the placement of the table here with the table in the printed output of this document.

Always use midrule to separate table header rows from data rows, and use it only for this purpose. This enables assistive technologies to recognise table headers and support their users in navigating tables more easily.

5 FIGURES

The “figure” environment should be used for figures. One or more images can be placed within a figure. If your figure contains third-party material, you must clearly identify it as such, as shown in the example below.

ACKNOWLEDGMENTS

To Alefiya, for the continual guidance and enthusiasm all throughout the project

REFERENCES

- [] Rafal Ablamowicz and Bertfried Fauser. 2007. *CLIFFORD: a Maple 11 Package for Clifford Algebra Computations, version 11*. Retrieved February 28, 2008 from <http://math.tntech.edu/rafal/cliff11/index.html>
- [] Patricia S. Abril and Robert Plant. 2007. The patent holder’s dilemma: Buy, sell, or troll? *Commun. ACM* 50, 1 (Jan. 2007), 36–44. <https://doi.org/10.1145/1188913.1188915>
- [] Sten Andler. 1979. Predicate Path expressions. In *Proceedings of the 6th. ACM SIGACT-SIGPLAN symposium on Principles of Programming Languages (POPL ’79)*. ACM Press, New York, NY, 226–236. <https://doi.org/10.1145/567752.567774>



Fig. 1. 1907 Franklin Model D roadster. Photograph by Harris & Ewing, Inc. [Public domain], via Wikimedia Commons. (<https://goo.gl/VLCRBB>).

- [] David A. Anisi. 2003. *Optimal Motion Control of a Ground Vehicle*. Master's thesis. Royal Institute of Technology (KTH), Stockholm, Sweden.
- [] Sam Anzaroot and Andrew McCallum. 2013. *UMass Citation Field Extraction Dataset*. Retrieved May 27, 2019 from <http://www.iesl.cs.umass.edu/data/data-umasscitationfield>
- [] Sam Anzaroot, Alexandre Passos, David Belanger, and Andrew McCallum. 2014. Learning Soft Linear Constraints with Application to Citation Field Extraction. arXiv:1403.1349
- [] Lutz Bornmann, K. Brad Wray, and Robin Haunschild. 2019. Citation concept analysis (CCA)—A new form of citation analysis revealing the usefulness of concepts for other researchers illustrated by two exemplary case studies including classic books by Thomas S. Kuhn and Karl R. Popper. arXiv:1905.12410 [cs.DL]
- [] Kenneth L. Clarkson. 1985. *Algorithms for Closest-Point Problems (Computational Geometry)*. Ph. D. Dissertation. Stanford University, Palo Alto, CA. UMI Order Number: AAT 8506171.
- [] Jacques Cohen (Ed.). 1996. Special issue: Digital Libraries. *Commun. ACM* 39, 11 (Nov. 1996).
- [] Sarah Cohen, Werner Nutt, and Yehoshua Sagie. 2007. Deciding equivalences among conjunctive aggregate queries. *J. ACM* 54, 2, Article 5 (April 2007), 50 pages. <https://doi.org/10.1145/1219092.1219093>
- [] Bruce P. Douglass, David Harel, and Mark B. Trakhtenbrot. 1998. Statecharts in use: structured analysis and object-orientation. In *Lectures on Embedded Systems*, Grzegorz Rozenberg and Frits W. Vaandrager (Eds.). Lecture Notes in Computer Science, Vol. 1494. Springer-Verlag, London, 368–394. https://doi.org/10.1007/3-540-65193-4_29
- [] Ian Editor (Ed.). 2007. *The title of book one* (1st. ed.). The name of the series one, Vol. 9. University of Chicago Press, Chicago. <https://doi.org/10.1007/3-540-09237-4>
- [] Ian Editor (Ed.). 2008. *The title of book two* (2nd. ed.). University of Chicago Press, Chicago, Chapter 100. <https://doi.org/10.1007/3-540-09237-4>

- [] Matthew Van Gundy, Davide Balzarotti, and Giovanni Vigna. 2007. Catch me, if you can: Evading network signatures with web-based polymorphic worms. In *Proceedings of the first USENIX workshop on Offensive Technologies (WOOT '07)*. USENIX Association, Berkeley, CA, Article 7, 9 pages.
- [] Torben Hagerup, Kurt Mehlhorn, and J. Ian Munro. 1993. Maintaining Discrete Probability Distributions Optimally. In *Proceedings of the 20th International Colloquium on Automata, Languages and Programming (Lecture Notes in Computer Science, Vol. 700)*. Springer-Verlag, Berlin, 253–264.
- [] David Harel. 1978. *LOGICS of Programs: AXIOMATICS and DESCRIPTIVE POWER*. MIT Research Lab Technical Report TR-200. Massachusetts Institute of Technology, Cambridge, MA.
- [] David Harel. 1979. *First-Order Dynamic Logic*. Lecture Notes in Computer Science, Vol. 68. Springer-Verlag, New York, NY. <https://doi.org/10.1007/3-540-09237-4>
- [] Lars Hörmander. 1985. *The analysis of linear partial differential operators. III*. Grundlehren der Mathematischen Wissenschaften [Fundamental Principles of Mathematical Sciences], Vol. 275. Springer-Verlag, Berlin, Germany. viii+525 pages. Pseudodifferential operators.
- [] Lars Hörmander. 1985. *The analysis of linear partial differential operators. IV*. Grundlehren der Mathematischen Wissenschaften [Fundamental Principles of Mathematical Sciences], Vol. 275. Springer-Verlag, Berlin, Germany. vii+352 pages. Fourier integral operators.
- [] IEEE. 2004. IEEE TCSC Executive Committee. In *Proceedings of the IEEE International Conference on Web Services (ICWS '04)*. IEEE Computer Society, Washington, DC, USA, 21–22. <https://doi.org/10.1109/ICWS.2004.64>
- [] Markus Kirschmer and John Voight. 2010. Algorithmic Enumeration of Ideal Classes for Quaternion Orders. *SIAM J. Comput.* 39, 5 (Jan. 2010), 1714–1747. <https://doi.org/10.1137/080734467>
- [] Donald E. Knuth. 1997. *The Art of Computer Programming, Vol. 1: Fundamental Algorithms (3rd. ed.)*. Addison Wesley Longman Publishing Co., Inc.
- [] David Kosiur. 2001. *Understanding Policy-Based Networking* (2nd. ed.). Wiley, New York, NY.
- [] Leslie Lamport. 1986. *TeX: A Document Preparation System*. Addison-Wesley, Reading, MA.
- [] Newton Lee. 2005. Interview with Bill Kinder: January 13, 2005. Video. *Comput. Entertain.* 3, 1, Article 4 (Jan.-March 2005). <https://doi.org/10.1145/1057270.1057278>
- [] Dave Novak. 2003. Solder man. Video. In *ACM SIGGRAPH 2003 Video Review on Animation theater Program: Part I - Vol. 145 (July 27–27, 2003)*. ACM Press, New York, NY, 4. <https://doi.org/99.9999/woot07-S422> <http://video.google.com/videoplay?docid=6528042696351994555>
- [] Barack Obama. 2008. A more perfect union. Video. Retrieved March 21, 2008 from <http://video.google.com/videoplay?docid=6528042696351994555>
- [] Poker-Edge.Com. 2006. Stats and Analysis. Retrieved June 7, 2006 from <http://www.poker-edge.com/stats.php>
- [] R Core Team. 2019. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria. <https://www.R-project.org/>
- [] Bernard Rous. 2008. The Enabling of Digital Libraries. *Digital Libraries* 12, 3, Article 5 (July 2008). To appear.
- [] Mehdi Saeedi, Morteza Saheb Zamani, and Mehdi Sedighi. 2010. A library-based synthesis methodology for reversible logic. *Microelectron. J.* 41, 4 (April 2010), 185–194.
- [] Mehdi Saeedi, Morteza Saheb Zamani, Mehdi Sedighi, and Zahra Sasanian. 2010. Synthesis of Reversible Circuit Using Cycle-Based Approach. *J. Emerg. Technol. Comput. Syst.* 6, 4 (Dec. 2010).
- [] Joseph Scientist. 2009. The fountain of youth. Patent No. 12345, Filed July 1st., 2008, Issued Aug. 9th., 2009.
- [] Stan W. Smith. 2010. An experiment in bibliographic mark-up: Parsing metadata for XML export. In *Proceedings of the 3rd. annual workshop on Librarians and Computers (LAC '10, Vol. 3)*, Reginald N. Smythe and Alexander Noble (Eds.). Paparazzi Press, Milan Italy, 422–431. <https://doi.org/99.9999/woot07-S422>
- [] Asad Z. Spector. 1990. Achieving application requirements. In *Distributed Systems* (2nd. ed.), Sape Mullender (Ed.). ACM Press, New York, NY, 19–33. <https://doi.org/10.1145/90417.90738>
- [] Harry Thornburg. 2001. *Introduction to Bayesian Statistics*. Retrieved March 2, 2005 from <http://ccrma.stanford.edu/~jos/bayes/bayes.html>
- [] TUG. 2017. *Institutional members of the TeX Users Group*. Retrieved May 27, 2017 from <http://www.tug.org/instmemb.html>
- [] Boris Veytsman. 2017. *acmart—Class for typesetting publications of ACM*. Retrieved May 27, 2017 from <http://www.ctan.org/pkg/acmart>