

Table 5. Summary statistics of the SMAPE values for /8, /16 and /24 networks (#net is the # of networks involved), where bold-font highlights the most accurate predictions among the different p -values with respect to a specific model.

Model	p -value	#net	Min	Q ₁	Median	Mean	SD	Q ₃	Max
/8 networks									
GC5VAR2	small	71	0.012	0.064	0.099	0.209	0.238	0.276	1.188
	medium	7	0.165	0.304	0.857	0.784	0.538	1.214	1.429
	large	3	0.104	0.553	1.003	0.845	0.676	1.216	1.429
GC4VAR2	small	71	0.001	0.109	0.202	0.300	0.315	0.368	2.000
	medium	12	0.145	0.269	0.478	0.628	0.451	0.891	1.429
	large	3	0.145	0.580	1.015	0.863	0.655	1.222	1.429
GC3VAR2	small	71	0.001	0.111	0.222	0.315	0.292	0.452	1.358
	medium	21	0.139	0.245	0.452	0.534	0.387	0.717	1.429
	large	9	0.126	0.389	0.474	0.797	0.584	1.429	1.714
GC2VAR2	small	70	0.060	0.134	0.265	0.458	0.507	0.590	2.000
	medium	53	0.067	0.170	0.383	0.719	0.690	0.898	2.000
	large	32	0.080	0.155	0.658	0.881	0.754	1.610	2.000
/16 networks									
GC5VAR2	small	7807	0.054	0.857	1.430	1.394	0.563	2.000	2.000
	medium	7808	0.043	0.759	1.310	1.262	0.598	1.910	2.000
	large	7765	0.021	0.791	1.340	1.272	0.596	2.000	2.000
GC4VAR2	small	7995	0.025	0.254	0.704	0.784	0.560	1.220	2.000
	medium	7840	0.009	0.637	1.140	1.201	0.602	1.710	2.000
	large	7803	0.021	0.645	1.140	1.206	0.602	1.710	2.000
GC3VAR2	small	6816	0.030	0.376	0.781	0.869	0.555	1.302	2.000
	medium	7897	0.009	0.571	1.140	1.129	0.601	1.710	2.000
	large	7849	0.011	0.571	1.140	1.136	0.602	1.710	2.000
GC2VAR2	small	7956	0.027	0.774	1.180	1.174	0.516	1.590	2.000
	medium	7903	0.022	0.780	1.200	1.196	0.530	1.640	2.000
	large	7887	0.011	0.770	1.180	1.192	0.532	1.640	2.000
/24 networks									
GC5VAR2	small	12408	0.006	0.857	1.290	1.243	0.572	1.710	2.000
	medium	5642	2.000	2.000	2.000	2.000	0.000	2.000	2.000
	large	5412	2.000	2.000	2.000	2.000	0.000	2.000	2.000
GC4VAR2	small	68	0.069	0.160	0.299	0.641	0.704	0.781	2.000
	medium	6810	2.000	2.000	2.000	2.000	0.000	2.000	2.000
	large	6619	2.000	2.000	2.000	2.000	0.000	2.000	2.000
GC3VAR2	small	10249	0.202	2.000	2.000	1.889	0.306	2.000	2.000
	medium	7735	2.000	2.000	2.000	2.000	0.000	2.000	2.000
	large	7578	2.000	2.000	2.000	2.000	0.000	2.000	2.000
GC2VAR2	small	12652	0.001	0.849	1.270	1.233	0.535	1.710	2.000
	medium	12583	0.005	0.857	1.270	1.232	0.539	1.710	2.000
	large	12549	0.001	0.857	1.270	1.233	0.542	1.710	2.000

Table 6. Summary statistics of SMAPE values of leveraging peer networks vs. sub-networks as helpers, where a indicates models leveraging sub-networks (i.e., time series data across multiple network resolutions), and bold-font highlights the more accurate prediction when comparing the prediction leveraging time series across network resolutions (i.e., leveraging sub-networks) and the prediction leveraging time series at a single network resolution (leveraging peer networks).

Model	Leveraging peer /8 networks vs. /16 sub-networks as helpers								Leveraging peer /16 networks vs. /24 sub-networks as helpers								
	net	Min	Q ₁	Q ₂	Mean	SD	Q ₃	Max	net	Min	Q ₁	Q ₂	Mean	SD	Q ₃	Max	
GCPenVAR	66	0.035	0.155	0.309	0.389	0.292	0.540	1.192	GCPenVAR	677	0.028	0.311	0.561	0.721	0.510	1.028	2.000
GCPenVAR a	66	0.026	0.154	0.273	0.388	0.301	0.525	1.209	GCPenVAR a	677	0.027	0.220	0.336	0.431	0.321	0.524	2.000
GCnVAR2	67	0.056	0.167	0.311	0.390	0.307	0.518	1.424	GCnVAR2	1350	0.006	0.334	0.609	0.796	0.555	1.185	2.000
GCnVAR2 a	67	0.069	0.397	0.589	0.748	0.527	1.009	1.980	GCnVAR2 a	1350	0.073	0.806	1.210	1.181	0.476	1.560	2.000
GCnVAR	54	0.065	0.402	2.000	1.373	0.829	2.000	2.000	GCnVAR	275	0.027	0.280	0.485	0.702	0.544	1.051	2.000
GCnVAR a	54	0.364	0.688	1.164	1.284	0.669	2.000	2.000	GCnVAR a	275	0.039	0.273	0.429	0.827	0.732	1.714	2.000
GCnVAR2	67	0.198	2.000	2.000	1.914	0.311	2.000	2.000	GCnVAR2	1292	0.167	1.143	1.714	1.557	0.548	2.000	2.000
GCnVAR2 a	67	0.391	2.000	2.000	1.946	0.249	2.000	2.000	GCnVAR2 a	1292	0.003	0.581	1.772	1.373	0.715	2.000	2.000
GC5VAR	39	0.077	0.155	0.284	0.396	0.307	0.558	1.110	GC5VAR	474	0.027	0.287	0.486	0.686	0.522	1.017	2.000
GC5VAR a	39	0.110	0.262	1.714	1.211	0.839	2.000	2.000	GC5VAR a	474	0.039	0.254	0.439	0.831	0.732	1.714	2.000
GC5VAR2	68	0.012	0.064	0.096	0.207	0.238	0.274	1.188	GC5VAR2	1302	0.088	1.140	1.710	1.501	0.579	2.000	2.000
GC5VAR2 a	68	0.068	0.196	0.401	0.791	0.751	1.571	2.000	GC5VAR2 a	1302	0.003	0.287	0.545	0.875	0.713	1.714	2.000
GC4VAR	54	0.017	0.134	0.254	0.562	0.647	0.641	2.000	GC4VAR	52	0.040	0.082	0.168	0.354	0.437	0.371	2.000
GC4VAR a	54	0.077	0.280	2.000	1.279	0.836	2.000	2.000	GC4VAR a	52	0.040	0.143	0.185	0.629	0.801	0.748	2.000
GC4VAR2	63	0.033	0.105	0.184	0.282	0.300	0.357	2.000	GC4VAR2	2	0.001	0.009	0.018	0.018	0.025	0.027	0.035
GC4VAR2 a	63	0.069	0.159	0.298	0.632	0.703	0.748	2.000	GC4VAR2 a	2	2.000	2.000	2.000	2.000	0.000	2.000	2.000
GC3VAR	45	0.061	0.136	0.219	0.357	0.290	0.491	1.111	GC3VAR	112	0.024	0.190	0.340	0.503	0.460	0.628	2.000
GC3VAR a	45	0.078	0.142	0.328	0.854	0.821	2.000	2.000	GC3VAR a	112	0.031	0.170	0.249	0.736	0.815	2.000	2.000
GC3VAR2	58	0.060	0.117	0.227	0.317	0.257	0.437	1.139	GC3VAR2	192	0.066	0.189	0.376	0.638	0.580	0.924	2.000
GC3VAR2 a	58	0.074	0.116	0.221	0.464	0.547	0.569	2.000	GC3VAR2 a	192	0.047	0.164	0.221	0.511	0.663	0.364	2.000
GC2VAR	47	0.032	0.132	0.209	0.318	0.265	0.392	1.058	GC2VAR	26	0.102	0.260	0.404	0.830	0.706	1.490	2.000
GC2VAR a	47	0.074	0.153	0.288	0.650	0.711	0.812	2.000	GC2VAR a	26	0.072	0.195	0.240	0.589	0.712	0.502	2.000
GC2VAR2	65	0.060	0.129	0.219	0.442	0.518	0.533	2.000	GC2VAR2	438	0.024	0.281	0.528	0.795	0.623	1.247	2.000
GC2VAR2 a	65	0.072	0.126	0.276	0.631	0.733	0.746	2.000	GC2VAR2 a	438	0.017	0.190	0.304	0.594	0.665	0.511	2.000

the state-of-the-art statistical techniques do not permit sound statistical models of sparse time series. (ii) We assume away the non-stationary time series, which is inherited from the notion of G-causality. Addressing these limitations is an important open problem. Our case study has two limitations that are imposed by the particular dataset. (iii) The dataset only lasts for 97 days. Although it is sufficient to demonstrate the usefulness of the framework, it would be better if we have access to a dataset with a longer period of time. (iv) The dataset is collected by low-interaction honeypots and therefore does not present rich semantic information about the attacks. Using datasets collected from high-interaction honeypots or production networks would resolve this issue, for which the framework is equally applicable.

6. Conclusion

We presented the CGC framework for characterizing the usefulness of G-causality in cybersecurity, especially its usefulness in predicting cyber attack rates. We investigated a number of models and drew a number of insights, which can be leveraged as a stepping-stone towards fully understanding the usefulness and limitations of G-causality in cybersecurity. There are many open problems for future research, including: How can we rigorously, rather than empirically, characterize the usefulness and limitations of G-causality? What are the utilities of other kinds of causality (e.g., Pearl-causality) in the cybersecurity domain? How can we model sparse and non-stationary time series in a principled fashion?

Acknowledgement. We thank the anonymous reviewers for their constructive comments. This work was supported in part by NSF Grant #1736209.

References

- [1] (2019), Internet security threat report. URL <https://www.symantec.com/content/dam/symantec/docs/reports/istr-24-2019-en.pdf>.
- [2] (2019), Personally identifiable information targeted in breaches that impact billions of records. URL <https://www.forgerock.com/resources/view/92170441/industry-brief/us-consumer-data-breach-report.pdf>.
- [3] ZHAN, Z., XU, M. and XU, S. (2013) Characterizing honeypot-captured cyber attacks: Statistical framework and case study. *IEEE Transactions on Information Forensics and Security* 8(11): 1775–1789.
- [4] ZHAN, Z., XU, M. and XU, S. (2014) A characterization of cybersecurity posture from network telescope data. In *International Conference on Trusted Systems* (Springer): 105–126.
- [5] ZHAN, Z., XU, M. and XU, S. (2015) Predicting cyber attack rates with extreme values. *IEEE Transactions on Information Forensics and Security* 10(8): 1666–1677.
- [6] CHEN, Y.Z., HUANG, Z.G., XU, S. and LAI, Y.C. (2015) Spatiotemporal patterns and predictability of cyberattacks. *PLoS one* 10(5): e0124472.
- [7] PENG, C., XU, M., XU, S. and HU, T. (2018) Modeling multivariate cybersecurity risks. *Journal of Applied Statistics* 0(0): 1–23.
- [8] BAKDASH, J.Z., HUTCHINSON, S., ZAROUKIAN, E.G., MARUSICH, L.R., THIRUMURUGANATHAN, S., SAMPLE, C., HOFFMAN, B. et al. (2018) Malware in the future? Forecasting of analyst detection of cyber events. *Journal of Cybersecurity* 4(1): ty007.
- [9] WERNER, G., YANG, S. and MCCONKY, K. (2017) Time series forecasting of cyber attack intensity. In *Proceedings of the 12th Annual Conference on cyber and information security research* (ACM): 18.
- [10] FANG, Z., ZHAO, P., XU, M., XU, S., HU, T. and FANG, X. Statistical modeling of computer malware propagation dynamics in cyberspace. *Journal of Applied Statistics*.
- [11] PRITOM, M., SCHWEITZER, K., BATEMAN, R., XU, M. and XU, S. (2020) Data-driven characterization and detection of covid-19 themed malicious websites. In *IEEE ISI'2020*.
- [12] PRITOM, M., SCHWEITZER, K., BATEMAN, R., XU, M. and XU, S. (2020) Characterizing the landscape of covid-19 themed cyberattacks and defenses. In *IEEE ISI'2020*.
- [13] FICKE, E. and XU, S. (2020) Apin: Automatic attack path identification in computer networks. In *IEEE ISI'2020*.
- [14] LI, D., LI, Q., YE, Y. and XU, S. (2020) Sok: Arms race in adversarial malware detection. *CoRR abs/2005.11671*.
- [15] XU, M., HUA, L. and XU, S. (2017) A vine copula model for predicting the effectiveness of cyber defense early-warning. *Technometrics* 59(4): 508–520.
- [16] FANG, Z., XU, M., XU, S. and HU, T. (2021) A framework for predicting data breach risk: Leveraging dependence to cope with sparsity. *IEEE Trans. Inf. Forensics Secur.* 16: 2186–2201.
- [17] LIU, Z., ZHENG, R., LU, W. and XU, S. (2021) Using event-based method to estimate cybersecurity equilibrium. *IEEE CAA J. Autom. Sinica* 8(2): 455–467.
- [18] GRANGER, C.W. (1969) Investigating causal relations by econometric models and cross-spectral methods. *Econometrica: Journal of the Econometric Society*: 424–438.
- [19] XU, L., ZHAN, Z., XU, S. and YE, K. (2014) An evasion and counter-evasion study in malicious websites detection. In *IEEE CNS*: 265–273.
- [20] LI, X., PARKER, P. and XU, S. (2011) A stochastic model for quantitative security analyses of networked systems. *IEEE Transactions on Dependable and Secure Computing* 8(1): 28–43.
- [21] XU, S., LI, X., PARKER, T. and WANG, X. (2011) Exploiting trust-based social networks for distributed protection of sensitive data. *IEEE T-IFS* 6(1): 39–52.
- [22] XU, L., ZHAN, Z., XU, S. and YE, K. (2013) Cross-layer detection of malicious websites. In *Third ACM Conference on Data and Application Security and Privacy (CODASPY'13)*: 141–152.
- [23] LI, Z., ZOU, D., XU, S., OU, X., JIN, H., WANG, S., DENG, Z. et al. (2018) Vuldeepecker: A deep learning-based system for vulnerability detection. In *Proc. NDSS'18*.
- [24] LI, Z., ZOU, D., XU, S., JIN, H., ZHU, Y., ZHANG, Z., CHEN, Z. et al. (2021), Vuldeelocator: A deep learning-based system for detecting and locating software

- vulnerabilities, IEEE TDSC (accepted for publication).
- [25] ZOU, D., WANG, S., XU, S., LI, Z. and JIN, H. (2019) μ vuldeepecker: A deep learning-based system for multiclass vulnerability detection. *IEEE Transactions on Dependable and Secure Computing* : 1–1doi:10.1109/TDSC.2019.2942930.
- [26] LI, Z., ZOU, D., XU, S., JIN, H., QI, H. and HU, J. (2016) Vulpecker: an automated vulnerability detection system based on code similarity analysis. In *Pro. ACSAC'16*: 201–213.
- [27] LI, Z., ZOU, D., XU, S., JIN, H., ZHU, Y., CHEN, Z., WANG, S. *et al.* (2021) Sysevr: A framework for using deep learning to detect software vulnerabilities. *IEEE Transactions on Dependable and Secure Computing (accepted for publication)* .
- [28] XU, S., YUNG, M. and WANG, J. (2021/04/28) Seeking foundations for the science of cyber security. *Information Systems Frontiers* doi:10.1007/s10796-021-10134-8.
- [29] XU, S. (2014) Cybersecurity dynamics. In *Proc. HotSoS'14*: 14:1–14:2.
- [30] XU, S. (2019) Cybersecurity dynamics: A foundation for the science of cybersecurity. In *Proactive and Dynamic Network Defense*, 1–31.
- [31] XU, S. (2020) The cybersecurity dynamics way of thinking and landscape (invited paper). In *ACM Workshop on Moving Target Defense*.
- [32] ZHENG, R., LU, W. and XU, S. (2018) Preventive and reactive cyber defense dynamics is globally stable. *IEEE TNSE* 5(2): 156–170.
- [33] LIN, Z., LU, W. and XU, S. (2019) Unified preventive and reactive cyber defense dynamics is still globally convergent. *IEEE/ACM Trans. Netw.* 27(3): 1098–1111.
- [34] HAN, Y., LU, W. and XU, S. (2020) Preventive and reactive cyber defense dynamics with ergodic time-dependent parameters is globally attractive. *CoRR* abs/2001.07958.
- [35] XU, S., LU, W. and ZHAN, Z. (2012) A stochastic model of multivirus dynamics. *IEEE Transactions on Dependable and Secure Computing* 9(1): 30–45.
- [36] XU, S., LU, W. and XU, L. (2012) Push- and pull-based epidemic spreading in networks: Thresholds and deeper insights. *ACM TAAS* 7(3).
- [37] XU, S., LU, W., XU, L. and ZHAN, Z. (2014) Adaptive epidemic dynamics in networks: Thresholds and control. *ACM TAAS* 8(4).
- [38] ZHENG, R., LU, W. and XU, S. (2015) Active cyber defense dynamics exhibiting rich phenomena. In *Proc. HotSoS*.
- [39] XU, M., SCHWEITZER, K.M., BATEMAN, R.M. and XU, S. (2018) Modeling and predicting cyber hacking breaches. *IEEE T-IFS* 13(11): 2856–2871.
- [40] KAR, M., NAZLIOĞLU, Ş. and AĞIR, H. (2011) Financial development and economic growth nexus in the mena countries: Bootstrap panel granger causality analysis. *Economic modelling* 28(1-2): 685–693.
- [41] SHOJAIE, A. and MICHAELIDIS, G. (2010) Discovering graphical granger causality using the truncating lasso penalty. *Bioinformatics* 26(18): i517–i523.
- [42] DEWAN, S. and RAMAPRASAD, J. (2014) Social media, traditional media, and music sales. *Mis Quarterly* 38(1).
- [43] TILGHMAN, P. and ROSENBLUTH, D. (2013) Inferring wireless communications links and network topology from externals using granger causality. In *MILCOM 2013-2013 IEEE Military Communications Conference (IEEE)*: 1284–1289.
- [44] DEKA, R.K., BHATTACHARYYA, D.K. and KALITA, J.K. (2019) Granger causality in tcp flooding attack. *IJ Network Security* 21(1): 30–39.
- [45] QIN, X. and LEE, W. (2004) Attack plan recognition and prediction using causal networks. In *20th Annual Computer Security Applications Conference (IEEE)*: 370–379.
- [46] ZAIONTZ, C. (2013) Real statistics using excel. cronbach's alpha. Retrieved January 21.
- [47] GRANGER, C.W. (1988) Some recent development in a concept of causality. *Journal of econometrics* 39(1-2): 199–211.
- [48] HIEMSTRA, C. and JONES, J.D. (1994) Testing for linear and nonlinear granger causality in the stock price-volume relation. *The Journal of Finance* 49(5): 1639–1664.
- [49] ASIMAKOPOULOS, I., AYLING, D. and MAHMOOD, W.M. (2000) Non-linear granger causality in the currency futures returns. *Economics Letters* 68(1): 25–30.
- [50] BROVELLI, A., DING, M., LEDBERG, A., CHEN, Y., NAKAMURA, R. and BRESSLER, S.L. (2004) Beta oscillations in a large-scale sensorimotor cortical network: directional influences revealed by granger causality. *Proceedings of the National Academy of Sciences* 101(26): 9849–9854.
- [51] YAMANISHI, K. and TAKEUCHI, J.I. (2002) A unifying framework for detecting outliers and change points from non-stationary time series data. In *Proceedings of the eighth ACM SIGKDD international conference on Knowledge discovery and data mining*: 676–681.
- [52] ROTHE, C. and SIBBERTSEN, P. (2006) Phillips-perron-type unit root tests in the nonlinear estar framework. *Allgemeines Statistisches Archiv* 90(3): 439–456.
- [53] CHEUNG, Y.W. and LAI, K.S. (1998) Power of the augmented dickey-fuller test with information-based lag selection. *Journal of Statistical Computation and Simulation* 60(1): 57–65.
- [54] LÜTKEPOHL, H. (2005) *New introduction to multiple time series analysis* (Springer Science & Business Media).
- [55] LIDDLE, A.R. (2007) Information criteria for astrophysical model selection. *Monthly Notices of the Royal Astronomical Society: Letters* 377(1): L74–L78.
- [56] AKAIKE, H. (1974) A new look at the statistical model identification. *IEEE transactions on automatic control* 19(6): 716–723.
- [57] LOMAX, R.G. (2007) *Statistical concepts: A second course* (Lawrence Erlbaum Associates Publishers).
- [58] CHEN, Z. and YANG, Y. (2004) Assessing forecast accuracy measures. *Preprint Series* 2010: 2004–10.
- [59] ALATA, E., DACIER, M., DESWARTE, Y., KAAÂNICHE, M., KORTCHINSKY, K., NICOMETTE, V., PHAM, V.H. *et al.* (2006) Collection and analysis of attack data based on honeypots deployed on the internet. In *Quality of Protection* (Springer), 79–91.
- [60] ALMOTAIRI, S., CLARK, A., MOHAY, G. and ZIMMERMANN, J. (2008) Characterization of attackers' activities in honeypot traffic using principal component analysis. In *2008 IFIP International Conference on Network and Parallel Computing (IEEE)*: 147–154.

- [61] NICHOLSON, W.B., MATTESON, D.S. and BIEN, J. (2017) VarX-l: Structured regularization for large vector autoregressions with exogenous variables. *International Journal of Forecasting* 33(3): 627–651.
- [62] GAO, Y., LI, Z. and CHEN, Y. (2006) A dos resilient flow-level intrusion detection approach for high-speed networks. In *ICDCS*, 6: 39.
- [63] DAGON, D., QIN, X., GU, G., LEE, W., GRIZZARD, J., LEVINE, J. and OWEN, H. (2004) Honeystat: Local worm detection using honeypots. In *International Workshop on Recent Advances in Intrusion Detection* (Springer): 39–58.
- [64] PHAM, V.H. and DACIER, M. (2011) Honeypot trace forensics: The observation viewpoint matters. *Future Generation Computer Systems* 27(5): 539–546.
- [65] ANTONATOS, S., POLAKIS, I., PETSAS, T. and MARKATOS, E.P. (2010) A systematic characterization of im threats using honeypots. In *ISOC Network and Distributed System Security Symposium (NDSS)*.
- [66] KREIBICH, C. and CROWCROFT, J. (2004) Honeycomb: creating intrusion detection signatures using honeypots. *ACM SIGCOMM computer communication review* 34(1): 51–56.
- [67] PORTOKALIDIS, G. and Bos, H. (2007) Sweetbait: Zero-hour worm detection and containment using low-and high-interaction honeypots. *Computer Networks* 51(5): 1256–1274.
- [68] ANAGNOSTAKIS, K.G., SIDIROGLOU, S., AKRITIDIS, P., XINIDIS, K., MARKATOS, E. and KEROMYTIS, A.D. (2005) Detecting targeted attacks using shadow honeypots .
- [69] PENG, C., XU, M., XU, S. and HU, T. (2017) Modeling and predicting extreme cyber attack rates via marked point processes. *Journal of Applied Statistics* 44(14): 2534–2563.
- [70] PROVOS, N. (2004) A virtual honeypot framework. In *Proc. USENIX Security Symposium*.
- [71] BALAS, E. and VIECCO, C.H. (2005) Towards a third generation data capture architecture for honeynets. *Proceedings from the Sixth Annual IEEE SMC Information Assurance Workshop* : 21–28.
- [72] INACIO, C. and TRAMMELL, B. (2010) Yaf: Yet another flowmeter. In *LISA*.