



Kandidatutkielma

Tietojenkäsittelytieteen kandiohjelma

NLP-hyökkäysten käyttökohteet

Akira Taguchi

21.2.2022

MATEMAATTIS-LUONNONTIETEELLINEN TIEDEKUNTA
HELSINGIN YLIOPISTO

Yhteystiedot

PL 68 (Pietari Kalmin katu 5)
00014 Helsingin yliopisto

Sähköpostiosoite: info@cs.helsinki.fi
URL: <http://www.cs.helsinki.fi/>

Sisältö

1	Sisällys	1
2	Historia	2
3	Hyökkäystaksonomia	3
4	Puolustusmetodit	4
5	Yhteenveto	5
	Lähteet	6
A	Sample Appendix	i

1 Sisällys

Ohjelmistojen hyökkäysrajapinta-ala kasvaa jatkuvasti. Osa haavoittuvaisuuksista korjataan heti havainnoinnin jälkeen, osa mitigoidaan ja osan vaikutusalue on manifestoituu vasta tulevaisuudessa . Luonnollisen kielen prosessointi (eng. Natural Language Processing, NLP) on osoittautunut hyväksi hyökkäysrajapinnaksi tätä teknologiaa hyödyntäviä osapuolia vastaan (Boucher et al., 2021). NLP-järjestelmät on tehty tulkitsemaan ihmisen luonnollista kieltä. Tämän kielen konekääntäminen aloitettiin jo vuonna 1949.

Tässä tutkielmassa tarkastellaan NLP-hyökkäysten käyttökohteita. Tähän kuuluu oleellisen historian esittely, hyökkäystaksonomia sekä puolustusmetodit. On tärkeää ymmärtää luonnollisen kielen prosessoinnin tarkoitus, jotta voidaan syventyä hyökkäyksiä mahdollistaviin ongelmiin sekä näiden ratkaisemiseen.

2 Historia

3 Hyökkäystaksonomia

4 Puolustusmetodit

5 Yhteenveto

Lähteet

Boucher, N., Shumailov, I., Anderson, R. ja Papernot, N. (2021). *Bad Characters: Imperceptible NLP Attacks*. arXiv: [2106.09898 \[cs.CL\]](#).

Liite A Sample Appendix

You can add one or more appendices to your thesis.