

Zloupotreba umjetne inteligencije

CERT.hr-PUBDOC-2023-12-410

Sadržaj

1	UVOD	3
2	MOGUĆNOSTI ZLOUPOTREBE	4
2.1	DEEPFAKE.....	4
2.1.1	<i>Opis</i>	4
2.1.2	<i>Zaštita</i>	5
2.2	AI GENERIRANI MEDIJSKI SADRŽAJ	6
2.2.1	<i>Opis</i>	6
2.2.2	<i>Zaštita</i>	7
2.3	AUTOMATIZIRANJE RAZVOJNOG CIKLUSA	8
2.3.1	<i>ChatGPT</i>	8
2.3.2	<i>WormGPT</i>	9
3	ZAKLJUČAK	10
4	LITERATURA	11

Ovaj dokument izradio je Laboratorij za sustave i signale Zavoda za elektroničke sustave i obradbu informacija Fakulteta elektrotehnike i računarstva Sveučilišta u Zagrebu.

Ovaj dokument vlasništvo je Nacionalnog CERT-a. Namijenjen je javnoj objavi te se svatko smije njime koristiti i na njega se pozivati, ali isključivo u izvornom obliku, bez izmjena, uz obvezno navođenje izvora podataka. Korištenje ovog dokumenta protivno gornjim navodima povreda je autorskih prava CARNET-a, a sve navedeno u skladu je sa zakonskim odredbama Republike Hrvatske.

1 Uvod

Umjetna inteligencija (*Artificial intelligence*, AI) (31) je sposobnost nekog sustava da oponaša ljudske aktivnosti kao što su zaključivanje, učenje, planiranje i kreativnost. Omogućuje tehničkim sustavima percipiranje okruženja i uzimanje u obzir onoga što vide te rješavanje problema.

Njezina je primjena vrlo široka u svakodnevnom životu. Primjerice, internetske trgovine koriste ju za pružanje personaliziranih preporuka svakom korisniku na temelju prethodnih pretraga, kupovina i mnogih drugih informacija o njegovom ponašanju na internetu. Isto tako, pretraživači interneta uče iz velike količine podataka o svakom korisniku kako bi pružili relevantne rezultate pretraživanja.

U Europi se radi i na naprednijim primjenama umjetne inteligencije kao što su autonomna vozila. Iako ona još nisu u široj upotrebi, novija vozila koriste sigurnosne funkcije na temelju umjetne inteligencije. Primjerice, vozača se upozori ako je tijekom vožnje napustio svoju traku, vozilo se automatski zaustavi ako se vozač pri parkiranju previše približi nekom objektu i slično. Također se istražuju mogućnosti primjene umjetne inteligencije u medicinskoj dijagnostici, optimiziranju proizvodnje u tvornicama, prepoznavanju kibernetičkih napada i slično.

Neočekivano brz razvoj umjetne inteligencije i njezina svestrana primjena doveli su do toga da su se razvile i neke nove zlonamjerne aktivnosti. Jedan od primjera je *deepfake* tehnologija koja manipulira video snimkama do te mjere da se savršeno realistično može prikazati događaj koji se nikad nije dogodio. Također, razvoj alata za automatiziranje razvijanja programskog kôda kao što je ChatGPT može se primijeniti u loše svrhe te automatizirati i ubrzati razvoj zlonamjernih programa.

2 Mogućnosti zloupotrebe

Postoje razni načini na koje se umjetna inteligencija može iskoristiti u zlonamjerne svrhe. Većina uspješnih napada se zasniva na generiranju ili manipulaciji sadržaja kojem je cilj na neki način obmanuti žrtvu. Osim toga, olakšan je i razvoj malicioznih programa, što potencijalno omogućuje zlonamjernim osobama slabijih tehničkih vještina lakše izvođenje napada. U nastavku su detaljnije opisane mogućnosti umjetne inteligencije, primjeri napada koje su pojedinci ili organizacije izvršili pomoću nje te razvoj kibernetičke sigurnosti u svrhu zaštite od takvih napada.

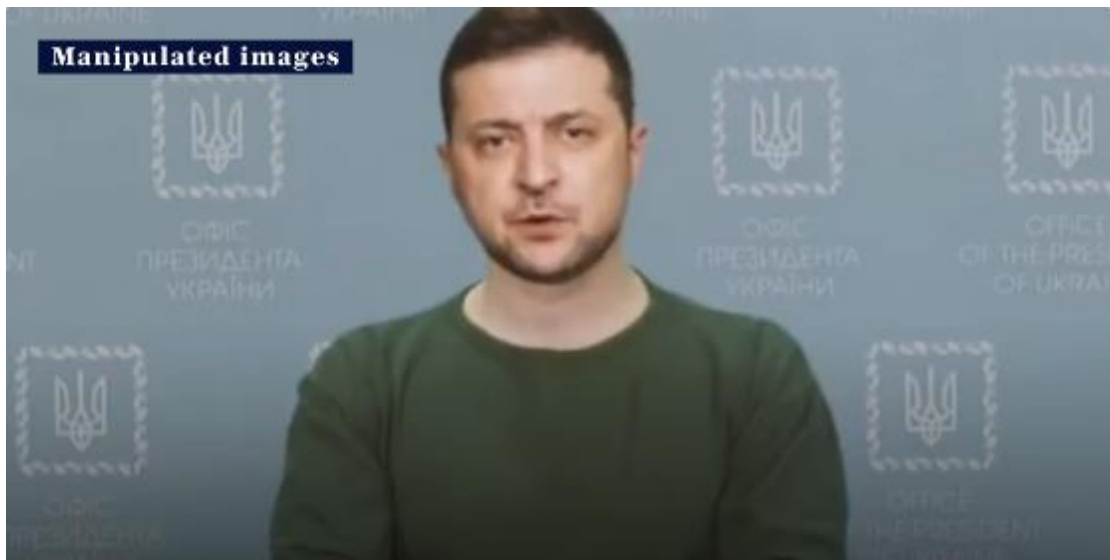
2.1 Deepfake

2.1.1 Opis

Deepfake jest fotografija, video ili audio sadržaj u kojemu je izgled/glas jedne osobe izmijenjen tako da izgleda/zvuči kao druga osoba (1) (2). Naziv dolazi od kombinacije riječi „deep learning“ i „fake“. Radi se o algoritmu strojnog učenja koji je u pravilu podijeljen na dva podalgoritma strojnog učenja (engl. machine learning, ML) modela. Ako se radi o izmjeni lica, prvi algoritam pokušava što bolje naučiti strukturu i karakteristike tog lica te kako ga najbolje preslikati na neko drugo u sklopu okruženja u kojem se osoba nalazi. Drugi je treniran da što bolje prepoznaje radi li se o falsifikatu ili ne. On uočava greške i nelogičnosti koje je prvi algoritam napravio te se time samostalno ispravlja i unapređuje. Trenutno postoji mnogo aplikacija poput Faceswapa (3) ili Faceit_livea (4) koje su besplatne. Da bi se koristile potrebno je imati grafičku karticu koja može podržati ML model.

Stvarni napadi

1. U neimenovanoj tvrtki krajem kolovoza 2019. godine, napadač je uspio uvjeriti izvršnog direktora da pošalje veliku svotu novca predstavljajući se kao njegov šef u telefonskom pozivu. Glas napadača bio je izmijenjen tako da je zvučao kao šef izvršnog direktora. Izgubljeno je 220 000€. (5)
2. FBI upozorava da je krađa osobnih informacija, u kombinaciji s deepfake tehnologijom, napadačima omogućila da uvjerljivo koriste identitet neke druge osobe, čime mogu doći na radna mjesta u organizaciji koja zapošljava na daljinu te ostvariti pristup povjerljivim informacijama. (6)
3. Krajem siječnja 2023. godine ruski hakeri dobili su pristup ukrajinskom web portalu i postavili deepfake video sadržaj predsjednika Zelenskog koji traži predaju Ukrajinaca. Video su ubrzo demantirali drugi mediji. (7)



Slika 1: Snimka zaslona deepfake videa predsjednika Zelenskog.
Tekst „Manipulated images“ je naknadno dodan ([izvor](#))

4. Sredinom 2022. godine pokrenut je phishing napad kasnije nazvan „ BitVex “. BitVex je i naziv platforme za kriptovalute putem koje su žrtve prevarene uplatom novca. Platforma je promovirana deepfake sadržajima koja je uključivala poznate ličnosti poput Elona Muska, Cathie Wood, Brada Garlinghousea, Michaela Saylora i Charlesa Hoskinsona. Procjena je da je ova prevara bila relativno neuspješna te da je izgubljeno samo 1700 američkih dolara. (8)



Slika 2: Snimka zaslona deepfake videa Elona Muska ([izvor](#))

2.1.2 Zaštita

Zaštita od ovakvih napada nije laka jer usavršavanjem tehnologije uvjerljivost manipuliranog sadržaja postaje sve veća.

Najčešći pokazatelji da se radi o deepfake-u su:

- okolina (npr. nepostojeće sjene, prejaki odsjaji) (9)
- nesavršenosti u licu (nerealistični madeži, nesinkronizirano treptanje, distorzije u unutrašnjosti usta poput nedostatka zubi i jezika, presavršeni zubi itd...) (10)
- nesinkroniziranost govora/zvuka i pomicanja usta npr. zbog kihanja

Nesinkroniziranost govora/zvuka i pomicanja usta može se gledati pri izricanju slova b, m i p. Ponekad nastaju pikseli sivih nijansi na rubovima izmijenjenih komponenti. Moguće je razaznati radi li se o krivotvorini i kada je osoba na snimci gledana iz drugog kuta. Ako pri izradi deepfake sadržaja nisu korištene fotografije osobe iz različitih kutova, algoritam ne može dokučiti izgled osobe iz drugog kuta te dolazi do distorzije. MIT (Massachusetts Institute of Technology) nudi [upitnik](#) za razlikovanje deepfake sadržaja od autentičnog u svrhu edukacije o ovoj temi.

U razvoju su i ML modeli koji raspoznaju autentičnost.

Klasični detektor deepfake vizualnog sadržaja temelji se na uočavanju grešaka koje su nastale manipuliranjem. Najčešće je to vezano uz analizu piksela koje ljudsko oko ne može vidjeti jer manipuliranjem slike rubni krajevi manipulirane komponente imaju posebne karakteristike. Jedan takav algoritam jest hibridni LSTM (Long short-term memory – neuronska mreža) i Encoder-Decoder algoritam (11). On radi na način da se paralelno: 1) analizira svaki pojedinačni piksel, 2) radi kompresija na razini cijele slike/videoa. Na kraju se uspoređuju rezultati obje funkcije i ako obje ukazuju na istu regiju materijal se smatra manipuliranim.

Postoje i detektori drukčijeg pristupa. Intelov detektor (12) se temelji na uočavanju uzoraka koje ljudsko oko ne vidi, a koji dokazuju valjanost sadržaja. Promatranjem suptilnih znakova poput širenja zjenica, promjena boja krvnih žila u skladu s otkucajima srca i sl. može se zaključiti je li sadržaj autentičan.

2.2 AI generirani medijski sadržaj

2.2.1 Opis

Osim mijenjanja medijskih sadržaja, moguće ih je i stvoriti, generirati. Veliki je broj područja na koja se može primijeniti AI generiranje, od jednostavnih tekstova, do slika pa čak i blogova ili novinskih članaka. Općenito, AI generirani medijski sadržaj generiran je NLP (Natural Language Processing) metodama čiji podskup čine NLU (Natural Language Understanding) i NLG (Natural Language Generation) metode (13). NLU pretvara upit korisnika iz tekstualnog oblika u strukturu podataka, a NLG iz neke strukture podataka u tekstualni oblik. Jedan od primjera tekstualnih generativnih modela jest ChatGPT. AI generativni algoritmi se dijele i na one koji koriste pristup Internetu zbog prikupljanja ključnog sadržaja te na one koji koriste samo internu obradbu.

Stvarni napadi

1. Početkom svibnja 2023. osoba s nadimkom „Mourningassassin“ putem platforme „Discord“ uspio je prodati AI generirane pjesme fanovima izvođača Franka Oceana za 13 000 američkih dolara. Pjesme su reklamirane kao neobjavljene pjesme koje su procurile. (14)
2. Krajem 2022. tvrtka „Graphika“ otkrila je operaciju „Spamouflage“, čiji je cilj bio širenje propagande kineske komunističke partije pomoću lažnog medija naziva „The Wolf News“. Generirano je svega 300 pregleda po videu. Kasnije je potvrđeno da je sadržaj generiran tehnologijom engleske AI tvrtke Synthesia. (15)



Slika 3: Snimka zaslona AI generiranog novinara.
Tekst „disinformation campaign video“ je naknadno dodan ([izvor](#))

3. Napadač šalje generirane e-mail poruke računovodstvu u kojima glumi zaposlenika tvrtke kako bi promijenio informacije o isplati plaće. U drugom primjeru slanjem e-mail poruke napadač pokušava dobiti korisničko ime i lozinku Facebook profila glumeći Facebook-ovu službenu korisničku podršku. (16)
4. Na društvenoj mreži X (prethodno Twitter) u svibnju 2023. otkrivena je botnet mreža nazvana Fox8. Mreža je bila sačinjena od 1140 lažnih profila. Koristila se za promoviranje kriptovaluta i dijeljenje zlonamjernih poveznica. Profili su bili vođeni chatbotom ChatGPT. Otkrivena je zbog popularne fraze „As an AI language model...“. Ovakva vrsta botneta predstavlja prijetnju za algoritam društvene mreže jer može zavarati procjenu zainteresiranosti za neku temu i time ju dodatno širiti. (17)

2.2.2 Zaštita

Razvijanje alata za otkrivanje generiranog medijskog sadržaja ima mnogo izazova i puno ga je teže razviti od algoritma za generiranje. U slučaju otkrivanja generiranih sadržaja,

često je riječ o raspoznavanju slika s interneta (18). To je zato što trenutno većina algoritama za razvijanje slika koriste sadržaje s Interneta koji se na kraju agregiraju u završnu sliku.

Generirani tekst je također teško raspoznati od autentičnog. Alati poput GPTZero ocjenjuju vjerojatnost da je tekst generiran pomoću dva kriterija: nedoumica i raspršenost (19).

Nedoumicom (engl. „perplexity“) algoritam mjeri koliko je upoznat s tekstom i koliko dobro može predvidjeti što slijedi nakon već zadanog niza riječi, tj. je li takav tekst bio dio njegovog skupa za treniranje.

Tekst s visokom nedoumicom, u odnosu na tekst s niskom, predstavlja niz riječi koje bi bile ocijenjene s malom vjerojatnošću pojave s obzirom na prethodni niz riječi. Modeli su općenito postavljeni tako da generiraju tekst s niskom nedoumicom.

Drugim kriterijem algoritam mjeri raspršenost (engl. „burstiness“) teksta, odnosno njegovu kaotičnost. Na primjer, tekst koji je napisao čovjek imaće razne i nepredvidive duljine rečenica, dok će u generiranom tekstu duljine rečenica biti predvidive i monotone.

Osim detekcije, pokušava se uvesti i normirano žigosanje (engl. „watermark“) generiranih sadržaja radi potiskivanja širenja dezinformacija. Google razvija SynthID, koji pri generiranju slika postavlja i digitalni vodeni žig (20). Taj žig je zapisan pikselima na način da je neprepoznatljiv ljudskom oku i ne može biti promijenjen filterima. Kad ga se pročita, ukazuje na to da je slika generirana.

2.3 Automatiziranje razvojnog ciklusa

Umjetnom inteligencijom zlonamjerne osobe nastoje automatizirati razvojni ciklus malicioznih programa. AI tekstualni generativni modeli poput ChatGPT-a su svojevrsni generatori programskog kôda koji se mogu upotrijebiti u te svrhe. Neki modeli su naknadno trenirani na skupu podataka malicioznog kôda, što ih čini još boljima u njihovom stvaranju.

2.3.1 ChatGPT

ChatGPT (21) je oblik generativnog AI-a, tzv. chatbot. Kreirala ga je tvrtka „OpenAI“ koju je 2015. godine osnovala grupa istraživača i poduzetnika poput Elona Muska i Sama Altmana. Jedan od investitora je i „Microsoft“.

ChatGPT i njegovi nasljednici rade tako da generiraju odgovor na upit korisnika služeći se velikim skupom podataka s interneta, na kojemu je model i treniran. Odgovor koji model generira je najučestaliji odgovor (po određenom kriteriju) iz tog skupa podataka.

2.3.1.1 Generiranje programa

Modelima poput ChatGPT-a ugrađene su sigurnosne mjere kako ne bi mogli odgovarati na upite određenih tema ili zahtjeva poput: pisanje malicioznog kôda, dijeljenje osobnih informacija, promoviranje nasilja, širenje dezinformacija itd.

Velika je prednost chatbota što korisnik ne mora poznavati neki od programskih jezika ili jezika za pretraživanje baze podataka, već se upiti mogu formirati prirodnim jezikom.

To omogućuje širu primjenu alata, no isto tako stvara i veći potencijal za zloupotrebu. Iz tog razloga ugrađene su sigurnosne mjere i zabranjene određene teme koje chatbot može posluživati. Korisnici često pokušavaju zaobići te sigurnosne mjere uz tzv. „jailbreak“. Ideja je da se specifičnim unosom pokušaju zaobići sigurnosne mjere koje omogućuju korištenje ChatGPT-a u zlonamjerne svrhe (npr. „grandma exploit“ (22)).

Tako je, na primjer, demonstriran razvoj malicioznog programa samo upitima ChatGPT-u (23). Pažljivim postavljanjem upita i rastavljanjem programa na više dijelova, demonstrator zaobilazi mjere zaštite chatbot-a i tako stvara maliciozni kôd.

2.3.1.2 Analiza sigurnosti kôda

Još jedna moguća upotreba ChatGPT-a je za analizu sigurnosti kôda. Cilj je otkriti ranjivosti kôda upitom ChatGPT-u.

U sljedećem ulomku opisan je testni primjer analize kôda koji se detaljnije može proučiti [ovdje](#). Isječci kôda su uzeti iz namjerno ranjive web aplikacije *Damn Vulnerable Web Application*.

U suštini, chatbot je uspješno otkrio ranjivosti generičkog tipa. Za kompliciranije probleme generirao je lažno pozitivne odgovore u koje je bio vrlo siguran, što je tipično za chatbotove, tzv. halucinacije (25). Zaključak je da nije proizveo uspješnije rezultate od postojećih alata za testiranje kôdova.

Dodatno je bio i tražen da ispravi dani kôd. Uspješno ga je modificirao i ispravio probleme generičkog tipa (jer oni često dolaze uz rješenja na Internetu) uz par iznimki.

2.3.2 WormGPT

WormGPT je bio chatbot baziran na jezičnom modelu otvorenog kôda GPT-J6B. Prvobitno je prodavan na „ HackForums “, forumu za kibernetičke alate i usluge, s rasponom cijena licence 500 – 5000 €. Kreirala ga je grupa od pet programera. Nije imao etičke mjere zaštite. Bio je sposoban napraviti maliciozne programe, lažne vijesti, a osobito je bio efikasan u kreiranju phishing napada uz dane informacije o kontekstu (26).



Slika 4: Snimka zaslona upita WormGPT-u ([izvor](#))

Razvojni programeri tvrde da su otkazali projekt zbog lošeg publiciteta.

3 Zaključak

Umjetna inteligencija je svojim naglim porastom već pokazala da otvara mnoga vrata kibernetičkom kriminalu. Trenutno su opasnosti uglavnom vezane uz socijalni inženjering, gdje se obmanjujućim sintetičkim sadržajem pokušava pojedinca ili širu publiku prevariti na neki način. Uz rastući problem sintetičkog sadržaja, porasla je i potreba za njegovim filtriranjem. Zato se razvijaju mnogi detektori s raznim metodologijama te se pokušava uvesti kultura umetanja digitalnog vodenog žiga pri generiranju takvog sadržaja. Uz napade, nastaju i chatbotovi koji imaju mogućnosti generiranja malicioznog kôda ili phishing e-mailova. Analiza kôda i modificiranje su u vrijeme pisanja dokumenta uspješni samo ako su problemi generičke vrste.

4 Literatura

1. [Mrežno] <https://www.aware.com/blog-how-does-deepfake-technology-work/> [Citirano 8.9.2023]
2. [Mrežno] <https://www.techslang.com/what-is-deepfake-technology/> [Citirano 8.9.2023]
3. [Mrežno] <https://faceswap.dev/> [Citirano 8.9.2023]
4. [Mrežno] https://github.com/alew3/faceit_live [Citirano 8.9.2023]
5. [Mrežno] <https://www.zdnet.com/article/forget-email-scammers-use-ceo-voice-deepfakes-to-con-workers-into-wiring-cash/> [Citirano 8.9.2023]
6. [Mrežno] <https://www.zdnet.com/article/fbi-warning-crooks-are-are-using-deepfakes-to-apply-for-remote-tech-jobs/> [Citirano 8.9.2023]
7. [Mrežno] <https://fortune.com/2023/04/14/deepfakes-ai-state-propaganda/> [Citirano 8.9.2023]
8. [Mrežno] <https://heimdalsecurity.com/blog/deep-fakes-of-elon-musk-promote-bitvex-fraud/> [Citirano 8.9.2023]
9. [Mrežno] <https://mindmatters.ai/2022/08/a-novel-trick-for-detecting-deepfakes-a-sideways-view/> [Citirano 8.9.2023]
10. [Mrežno] <https://www.media.mit.edu/projects/detect-fakes/overview/> [Citirano 8.9.2023]
11. [Mrežno] <https://ieeexplore.ieee.org/stamp/stamp.jsp?arnumber=8626149&tag=1> [Citirano 8.9.2023]
12. [Mrežno] <https://www.intel.com/content/www/us/en/newsroom/news/intel-introduces-real-time-deepfake-detector.html#gs.ie8qbh> [Citirano 8.9.2023]
13. [Mrežno] <https://www.xenonstack.com/blog/difference-between-nlp-nlu-nlg> [Citirano 8.9.2023]
14. [Mrežno] <https://www.vice.com/en/article/z3mn75/scammer-made-thousands-selling-leaked-frank-ocean-tracks-that-were-fake-ai-generated-the-line-steer-it> [Citirano 8.9.2023]
15. [Mrežno] <https://therecord.media/deepfake-news-anchors-spread-chinese-propaganda-on-social-media> [Citirano 8.9.2023]
16. [Mrežno] <https://abnormalsecurity.com/blog/generative-ai-chatgpt-enables-threat-actors-more-attacks> [Citirano 8.9.2023]
17. [Mrežno] https://arstechnica.com/information-technology/2023/08/chatgpt-boosts-crypto-botnet-with-ai-generated-tweets/?utm_brand=arstechnica&utm_source=twitter&utm_social-type=owned&utm_medium=social [Citirano 8.9.2023]
18. [Mrežno] <https://chatonai.org/ai-image-detection-tools> [Citirano 8.9.2023]
19. [Mrežno] <https://zapier.com/blog/ai-content-detector/> [Citirano 8.9.2023]
20. [Mrežno] <https://edition.cnn.com/2023/08/30/tech/google-ai-images-watermark/index.html> [Citirano 8.9.2023]
21. [Mrežno] <https://www.techtarget.com/whatis/definition/ChatGPT> [Citirano 28.9.2023]
22. [Mrežno] <https://kotaku.com/chatgpt-ai-discord-clyde-chatbot-exploit-jailbreak-1850352678> [Citirano 28.9.2023]
23. [Mrežno] <https://www.forcepoint.com/blog/x-labs/zero-day-exfiltration-using-chatgpt-prompts> [Citirano 28.9.2023]
24. [Mrežno] <https://www.edureka.com/blog/steganography-tutorial> [Citirano 28.9.2023]
25. [Mrežno] <https://edition.cnn.com/2023/08/29/tech/ai-chatbot-hallucinations/index.html> [Citirano 28.9.2023]

26. [Mrežno] <https://krebsonsecurity.com/2023/08/meet-the-brains-behind-the-malware-friendly-ai-chat-service-wormgpt/> [Citirano 28.9.2023]
27. [Mrežno] <https://mindmatters.ai/2022/08/a-novel-trick-for-detecting-deepfakes-a-sideways-view/> [Citirano 28.9.2023]
28. [Mrežno] <https://www.media.mit.edu/projects/detect-fakes/overview/> [Citirano 28.9.2023]
29. [Mrežno] <https://research.nccgroup.com/2023/02/09/security-code-review-with-chatgpt/> [Citirano 28.9.2023]
30. [Mrežno] <https://www.cp24.com/news/trudeau-said-that-he-invested-in-the-same-thing-how-a-deepfake-video-cost-an-ontario-man-11k-us-1.6559497> [Citirano 28.9.2023]
31. [Mrežno]
https://www.europarl.europa.eu/pdfs/news/expert/2020/9/story/20200827STO85804/20200827STO85804_hr.pdf [Citirano 9.11.2023]