

## CSC:n asiantuntijalausunto tulevaisuusvaliokunnalle (E 27/2024 vp)

Näkökulmia supertietokoneilla tehtävän tutkimuksen hyödyntämiseen EU:n terveysunionissa pitkällä aikavälillä

CSC – Tieteen tietotekniikan keskus Oy (CSC) kiittää mahdollisuudesta lausua suurteholaskennan eli supertietokoneilla tehtävän tutkimuksen hyödyntämisestä terveysunionissa. CSC pitää tervetulleina Suomen ennakkovaikuttamista Euroopan terveysunionin tulevaisuuteen sekä Suomen kantoja mm. keskittymisestä terveysalalla EU-sääntelyn laatuun ja sääntelyn ja muiden tavoitteiden toimeenpanoon, EU:n terveystietoavaruuden (EHDS) alaisten tietojen toisiokäytön mahdollistamiseen TKI-toiminnan hyväksi, aktiiviseen osallistumiseen miljoonan genomien hankkeeseen sekä ylipäänsä kokonaisvaltaiseen lähestymiseen terveysturvallisuuden, tautien hoidon ja ennaltaehkäisyn ja esimerkiksi lääkkeiden kehittämisen osalta.

### Terveyksdatan hyödyntäminen TKI-toiminnassa edellyttää mahdollistavaa lainsäädäntöä ja rakenteellisia uudistuksia

Suomen kilpailukyvyllä on tärkeää, että suurteholaskentaa voidaan kattavasti hyödyntää tutkimuksessa terveysalan kasvattamiseksi. Kansallisen datan ja laskennan infrastruktuurin, EuroHPC-supertietokone LUMIn sekä sitä seuraavien suurteholaskennan investointien ansiosta Suomella onkin käytössään maailman kärkeä olevat TKI-työkalut terveysunionin ja esimerkiksi kansallisen terveys- ja hyvinvointialan kasvua ja uudistumista vauhdittavan TKI-ohjelman tavoitteiden edistämiseen. Tämän potentiaalin hyödyntämiseksi lainsäädännön ja määräysten on luotava mahdollisuuksia terveystiedon tietoturvalle käytölle terveysalan TKI-toiminnassa, joka on yhä kansainvälisempää ja yhä dataintensiivisempää.

Yhtäläisten osallistumismahdollisuuksien turvaamiseksi ja tutkimuksen edistämiseksi EHDS:ssä, sosiaali- ja terveystietojen toissijaista käyttöä säätelevää kansallista lakia on uudistettava siten, että aineistojen ja datan uudelleenkäyttö ja yhdistely helpottuvat. Toisiolain nojalla annettava määrästä tietoturvalisistä käyttöympäristöistä tulee päivittää, jotta se sallii kansallisesti tehtyjen suurteholaskennan investointien käytön sensitiivisellä datalla, kuten terveysdatalla tehtävään tutkimukseen. Tämän rinnalla on jatkettava EU-tason vaikuttamista sen varmistamiseksi, että terveystietoavaruutta koskevan asetuksen nojalla asetettavat EU-tason vaatimukset turvallisille käyttöympäristöille toimivat samansuuntaisesti. Lisäksi on tärkeää varmistaa, että terveystiedon TKI-käytön näkökulma huomioidaan EU:n yleisen tietosuojasetuksen väliarviossa sekä hallitusohjelmassa linjatussa kansallisen tietosuojalainsäädännön uudistuksessa, sillä tietosuojasääntely asettaa olennaisia reunaehdoja terveystiedon hyödyntämiselle.



## **Terveysturvallisuudesta sairauksien parempaan ymmärtämiseen, diagnosointiin ja hoitoon – supertietokoneet unionin tavoitteiden mahdollistajina**

Viimeistään COVID-19 -pandemia paljasti, että EU:n tavoitteet strategisesta autonomiasta sekä kansalaistensa hyvinvoinnin, terveyden ja turvallisuuden edistämisestä kytkeytyvät yhä enemmän biologisten ja lääke- ja terveystieteellisten alojen tutkimus-, kehitys- ja innovaatiotoimintaan sekä sitä mahdollistaviin data- ja tutkimusinfrastruktuureihin ja digitaalisiin teknologioihin. Esimerkiksi CSC:n suurteholaskentaresurssija on hyödynnetty viruksen ilmaitse tapahtuvan leviämisen mallintamiseen ja toimintamekanismin ymmärtämiseen suojautumistoimien parantamiseksi ja lääketutkimuksen pohjustamiseksi<sup>1</sup>. Eurooppalaiseen dataportaaliiin<sup>2</sup> on kerätty ainakin kahdeksan miljoonaa koronavirussekvenssiä, minkä ansiosta infektion molekyyli-tason toiminnan ja virusvarianttien tutkimusta voidaan jatkaa tulevaisuuden varautumistarpeita varten.

Supertietokoneiden ympärille kehittyvät, esimerkiksi yhteiseurooppalaisten EuroHPC-supertietokonekeskusten kaltaiset ekosysteemit tuovat etua monitieteiselle tutkimukselle mahdollistaessaan eri alojen suurten data-aineistojen yhdistelyn ja käsittelyn, mikä luo mahdollisuuksia myös E-kirjeessä korostetulle Yhteinen terveys (*one health*) -lähestymistavalle. Kajaanissa sijaitsevalla EuroHPC-supertietokone LUMilla on jo käynnissä merkittäviä biodiversiteetti-, ympäristö- ja ilmastodataan perustuvia tutkimushankkeita<sup>3, 4</sup>, ja yhdistelemällä tällaista tietoa terveys- ja lääketieteelliseen tietoon voidaan tulevaisuudessa oppia enemmän esimerkiksi lääkeresistenssin ja pandemioiden kaltaisista terveysuhista. Ilmastopuolella superkoneilla on jo tehty esimerkiksi sään ääri-ilmiöiden mallintamista<sup>5</sup> ja niitä voidaan tulevaisuudessa hyödyntää myös rajat ylittävien terveysuhkien, kuten vaikkapa säteilytilanteiden, yhä tarkempaan mallintamiseen.

## **Tarkempi sairauksien analyysi edellyttää monipuolisempaa dataa ja suurempaa laskentatehoa**

Eri alojen suuria data-aineistoja yhdistelemällä on siis mahdollista saavuttaa merkittäviä tieteellisiä läpimurtoja esimerkiksi juuri terveystutkimuksessa. Dataintensiivinen tutkimus on yleistynyt jonkin aikaa, ja nykyään on tyypillistä myös perinteiseen rekisteritutkimukseen perustuvien havaintojen täydentäminen eri biologisella tiedolla, kuten proteiini-, genomi- tai muulla molekyyli-tason datalla. Tällaisella analyysillä voidaan yhä tarkemmin tunnistaa harvinaissairauksia ja sairauksien syntymekanismia sekä kehittää yksilöllistetympiä hoitoja ja lääkityksiä, mikä kaikki tukee terveystietoyhteistyön tavoitteita sairauksien ehkäisyn ja hoidon vahvistamisesta.

<sup>1</sup> <https://csc.fi/uutinen/puhdilla-koronaa-vastaan-katsaus-covid-19-pikakaistan-tuloksiin/> ; <https://www.aka.fi/en/about-us/whats-new/press-releases/20202/high-level-research-infrastructures-support-covid-19-research/>

<sup>2</sup> <https://www.covid19dataportal.org/the-european-covid-19-data-platform>

<sup>3</sup> <https://biodt.eu/>

<sup>4</sup> <https://stories.ecmwf.int/destination-earth/index.html>

<sup>5</sup> <https://stories.ecmwf.int/m-t-o-france-wins-bid-to-develop-destination-earth-s-on-demand-extremes-digital-twin/index.html>



Terveytutkimuksessa käytettävien aineistojen monipuolistuminen ja monimutkaistuminen edellyttävät paitsi saatavilla olevan datan määrän ja laadun parantamista, myös sellaista tallennus- ja laskentakapasiteettia ja erityisosaamista, jota vain maailman tehokkaimmat supertietokonekeskukset, kuten LUMI ja sen ympärille rakentuva ekosysteemi, kykenevät tarjoamaan. Yhteiskunnallisesti vaikuttavan terveytutkimuksen edistämisen kannalta onkin hyvin positiivista, että Suomi tavoittelee terveysunionissa EHDS-toisiokäytön edistämistä sekä aktiivisuutta miljoonan genomien aloitteen toimeenpanossa ja on lisäksi jo nyt sitoutunut uuteen EuroHPC-supertietokonehankintaan nykyisen LUMI:n tullessa käyttökänsä päähän muutaman vuoden kuluttua.

### **Supertietokoneet ovat tekoälykehityksen kasvun edellytys myös tulevaisuudessa**

Tarve suurteholaskennalle kasvaa koko ajan, kun tekoäly ja mutkikkaat matemaattiset mallit yleistyvät. Supertietokoneilla on erityisen keskeinen rooli tekoälyn kehittämisessä, joka perustuu viime kädessä suuriin ja laadukkaisiin data-aineistoihin. Uusien supertietokoneiden suuri laskentateho tekoälysovelluksissa nojaa GPU-prosessoreihin perinteisten CPU-prosessoreihin sijaan ja niillä on saatu lyhennettyä laskentaprosesseja alle kymmenykseen aiemmasta esim. kuvaperusteisten, syöpätutkimusta tukevien algoritmien kehityksessä. Tällä tavalla LUMI-supertietokone on edistänyt esimerkiksi syöpädiagnostiikan tehostamista<sup>6,7</sup>, mikä on lupaavaa EU:n syöväntorjuntasuunnitelmalle osana terveysunionia myös tulevaisuudessa tekoälytutkimuksen kehittyessä edelleen.

Tekoälyä voi hyödyntää hyvin laajasti terveysunionin tavoitteissa, ja esimerkiksi kielimalleihin perustuvan generatiivisen tekoälyn käyttö terveysalalla tulee hyödyttämään sekä hoito- että TKI-työtä. Sen avulla voidaan esimerkiksi nopeuttaa tiedonkeruuta hoitosuunnitelmien seurannassa tai tutkimuslähteiden tarkastamisessa.<sup>8</sup> Eurooppalaiselle terveysunionille on siten tärkeää, että kielimalleja kehitetään tulevaisuudessa yhä vahvemmin eurooppalaisilla kielillä alaa dominoivan englannin rinnalla. LUMIlla on jo kehitetty eurooppalaisia kielimalleja pohjatyöksi niiden laajentamista ja alakohtaista jatkojalostusta varten<sup>9</sup>, mukaan lukien suomeksi<sup>10</sup> ja muilla pohjoismaisilla kielillä<sup>11</sup>.

Yllä mainittu datalähteiden monipuolistuminen mm. EHDS:n ja miljoonan genomien hankkeen onnistuneen toimeenpanon seurauksena yhdistettynä suureen laskentatehoon luovat merkittävää lisäpotentiaalia, sillä generatiivisen tekoälyn kunnollinen hyödyntäminen hoito- ja diagnosointityön tukena vaatii yksityis- ja myös yksilökohtaisempaa dataa.<sup>12</sup> Samalla on huomattava, että inhimillinen

<sup>6</sup> <https://www.lumi-supercomputer.eu/lumi-provides-new-opportunities-for-artificial-intelligence-research/>

<sup>7</sup> <https://www.hpcwire.com/2023/06/08/inside-comptai-one-of-lumis-first-projects/>

<sup>8</sup> <https://www.brookings.edu/articles/generative-ai-in-health-care-opportunities-challenges-and-policy/>

<sup>9</sup> <https://www.lumi-supercomputer.eu/now-language-models-will-be-trained-in-european-languages-on-one-of-the-worlds-largest-text-collections/>

<sup>10</sup> <https://www.lumi-supercomputer.eu/research-group-created-the-largest-finnish-language-model-ever-with-the-lumi-supercomputer/>

<sup>11</sup> <https://www.silo.ai/blog/viking-7b-13b-33b-sailing-the-nordic-seas-of-multilinguality>

<sup>12</sup> Ks. viite 9.



osaaminen ja ymmärrys tekoälyn kehittämisestä ja hyödyntämisestä tutkimuksessa ja hoitotyössä vaativat panostuksia myös pelkän ”raudan”, datan ja teknologian lisäksi.

## **Panostukset tietoturvalliseen datanhallintaan ja mahdollistavaan sääntelyyn maksimoivat supertietokoneiden tuoman kilpailuedun terveystieteen TKI-toiminnalle**

Koska hyvin hallittu data on niin suurteholaskennan kuin tekoälyn kehittämisenkin raaka-ainetta, Suomen ja EU:n tulee panostaa datan saatavuuteen, uudelleenkäytettävyyteen ja hallintaan. *Lainsäädännön tulee tukea ja mahdollistaa näitä tavoitteita, jotta tutkimusinfrastruktuuri-investoinneista ja työkaluista, kuten supertietokoneista, saadaan kaikki hyöty irti terveystieteen TKI-toiminnassa. Tulee varmistaa, että eurooppalaiset ja kansalliset vaatimukset ja määräykset datan käsittely-ympäristöistä, jollainen supertietokonekin on, tukevat sekä tietoturvan että käytettävyyden tavoitteita. Niiden tulee tulevaisuudessa mahdollistaa suurteholaskennan käyttö kattavasti ja monipuolisesti myös sellaisiin terveystietoihin, jotka ovat tietosuojasääntelyn ja toisilain alaisia.* COVID-19 tuskin jää viimeiseksi pandemiaksi, ja EuroHPC-järjestelmät voivat olla keskeisessä roolissa sekä itse sairauksien tutkimisessa että terveydenhuollon resurssien optimointia koskevien laskelmien tekemisessä, mikäli ne saadaan osaksi tietoturvallisen datanhallinnan järjestelmää. Tämä tarkoittaa kansalaisista kerätyn terveystiedon laaja-alaista hyödyntämistä tavalla, jota ei välttämättä voida toteuttaa hyvinvointialueiden tai sairaaloiden omissa järjestelmissä. Myös tekoälyn vaikuttava hyödyntäminen hoitotyössä suomen kielellä tulee edellyttämään suurteholaskennan sallimista nykyistä monipuolisemmille datalähteille.

Supertietokoneet ovat keskitettyjä mutta monikäyttöisiä, jaettuina järjestelmiä, minkä takia vaaditaan mekanismeja suojaamaan arkaluonteista henkilötietoa muilta käyttäjiltä ja ulkopuolisilta. *Moni uusi tietoturvamekanismi on vasta kehitteillä, joten määräysten tulee asettaa ja kuvailla tietoturvan tila ja tavoitteet määräämättä kuitenkaan yksityiskohtaisesti keinoja, joilla ne tulee saavuttaa. Esimerkiksi nykyisten, toisilain nojalla annettavien määräysten kategorinen kieltö internetyhteyksistä henkilötietoon perustuvan terveystiedon käsittelyssä sulkee pois monta potentiaalista salausmenetelmää sekä käytännössä nykyisten supertietokoneressurssien tehon hyödyntämisen toisilain alaisessa tutkimuksessa.*

*Teknologianeutraalius on ollut EU-sääntelyssä vallitseva periaate, ja tätä ajatusta tulee noudattaa myös kansallisella tasolla. Periaatteellisella tasolla on pohdittava, sopiiko salaisten aineistojen suojaamiseen tehty Katakri-työkalu täysin lähtökohdaksi tieteelliseen tutkimuskäyttöön liittyvien järjestelmien arviointiin. Hallitusohjelman mukaiset tietosuojalainsäädännön ja toisilain uudistukset, EHDS sekä käsittely-ympäristöjen vaatimusten päivitys tarjoavat mahdollisuuden ottaa tietoturvalliseen käyttöön suurteholaskennan tutkimusresurssit, joiden merkitys suomalaisen ja eurooppalaisen terveystieteen tutkimukselle ja tulevaisuudelle on kiistaton.*

Lainsäädännön tarkastelun ja esteiden purkamisen lisäksi on tärkeää kehittää entistä toimivampia yhteyksiä eurooppalaisten datainfrastruktuurien ja EuroHPC-supertietokonekeskusten välillä sekä hajautetun arkkitehtuurin avulla että tavoittelemalla datainfrastruktuurien sijoittautumista



supertietokoneiden yhteyteen sekä yleisesti että bio- ja terveystieteissä<sup>13</sup>. Tietoturvallisen datanhallinnan ja laskennan kehitystarpeet tutkimuskäyttöä varten tulevat painottumaan EHDS:n myötä, ja tällainen työ on Suomessa jo käynnissä<sup>14</sup>. Suomi on jo nyt edelläkävijä datanhallintainfrastruktuurien ja supertietokoneiden yhdistämisessä kansallisen datanhallinnan ja laskennan yhteishankkeen myötä sekä operoidessaan bioinformatiikkaan erikoistunutta Suomen ELIXIR-osakeskusta<sup>15</sup>. Yleisesti ottaen suurteholaskennan käyttö tietosuojaan tai aineettomien oikeuksien alaisille aineistoille esimerkiksi ennustavien, ml. tekoälymallien tapauksessa edellyttää, että palvelukehityksessä on ison laskentakapasiteetin lisäksi käyttötarkoitusta huomioivaa datanhallinnan ja tietosuojaan osaamista.

Espoossa, 21.5.2024

CSC-Tieteen tietotekniikan keskus Oy

Kimmo Koski

Toimitusjohtaja

Irina Kupiainen

Yhteiskuntasuhdejohtaja

---

<sup>13</sup> Vrt. <https://csc.fi/uutinen/cscn-raskaan-laskennan-ymparisto-laajenee-massiivisella-datanhallintajarjestelmalla/>

<sup>14</sup> <https://csc.fi/blogi/parempia-yksilollisia-hoitomenetelmia-kehittyneiden-tiedonsiirto-ja-analyysimenetelmien-avulla/>

<sup>15</sup> <https://elixir-europe.org/>; Valtiosopimus 7/2015

