

Tweb id:

**2nd course on bioinformatics tools for Next Generation Sequencing data mining:
use of bioinformatics tools for typing pathogenic *E. coli***

Aika: 16.–17.6.2016
Järjestäjä & paikka: EU Reference Laboratory for *E. coli* (EU-RL VTEC)
Istituto Superiore di Sanità, Rooma, Italia
Verkkosivu: <http://www.iss.it/vtec/index.php?lang=2&id=248&tipo=22>
(kurssimateriaalit tulevat tänne)

EURL VTEC järjesti kansallisille vertailulaboratorioille suunnatun kurssin, jossa esiteltiin erityisesti STEC-bakteerien kokogenomisekvenssidatan analysointiin soveltuvia, selainpohjaisia työkaluja. Osa työkaluista oli käytettävissä EURL:n ylläpitämällä Galaxy-pohjaisella ARIES-alustalla (<https://w3.iss.it/site/aries/>) ja osa Tanskan teknillisen yliopiston (DTU) genomiepidemiologiakeskuksen (CGE) ylläpitämällä sivustolla (<http://www.genomicepidemiology.org/>). Kummatkin palvelut ovat ilmaisia, mutta ARIES vaatii rekisteröitymisen. Kurssi koostui luennoista ja tietokoneharjoituksista seuraavilla aihealueilla:

1. EFSA:n ja ECDC:n kokogenomisekvenssointiin (WGS) liittyvät aktiviteetit
2. Galaxy-arkkitehtuuri ja ARIES-klusteri
3. Virulotyyppitys
4. Serotyyppitys
5. wgSNP-tyypitys (whole genome single nucleotide polymorphism)
6. HReVAP (high resolution virulence allelic profiling)
7. MLST ja agMLST (accessory genome multi-locus sequence typing)

Aluksi **EFSA:n ja ECDC:n** edustajat esittelivät WGS-aktiviteettejaan: järjestettyjä kokouksia, julkaistuja asiantuntijalausuntoja ja rahoitettuja hankkeita. ECDC:llä FWD-Next asiantuntijaryhmä on julkaissut asiantuntijalausunnon aiheesta (<http://ecdc.europa.eu/en/publications/Publications/food-and-waterborne-diseases-next-generation-typing-methods.pdf>) ja julkaisee pian ECDC:n WGS-strategian epidemianselvityksiin. Haasteena on vertailukelpoisen tyyppitysdatan saaminen tulevaisuudessa, kun osa kansallisista laboratorioista siirtyy WGS-menetelmiin, kun taas osa pitäytyy vanhoissa menetelmissä. ECDC oli teettänyt kyselyn WGS-menetelmiin siirtymisestä kansallisille laboratorioille ja EFSA suunnitteli toteuttavansa samanlaisen kyselyn EURL:n avustuksella. Tarvittaessa ECDC teettää sekvenssointeja kansallisten laboratorioiden lähettämistä kannoista Source Bioscience (UK) -yrityksellä. Toistaiseksi ECDC:llä ei ole kuitenkaan kerätty tai analysoitu STEC-bakteerien WGS-dataa.

Galaxy on avoimen lähdekoodin IT-arkkitehtuuri, jota EURL käyttää **ARIES**-alustassaan. ARIES tarjoaa selainpohjaisen, graafisen käyttöliittymän ja erityisesti STEC-bakteerien WGS-datan analysoimiseen tarkoitettuja työkaluja. ARIES:in tarjoamia työkaluja käyttäjä voi itse automatisoida työvuoksi (pipeline/workflow) visuaali-

sella työkalulla (Workflow canvas). Galaxy-arkkitehtuuri mahdollistaa myös kotitekoisten työkalujen kehityksen, joiden lisäämistä ARIES:iin voi pyytää EURL:ltä. EURL ei kerää eikä varmuuskopioi käyttäjien ARIES:iin lataamaa dataa. ARIES on ollut julkisessa käytössä syyskuusta 2015 lähtien ja sillä on tällä hetkellä 40 käyttäjää (3 Suomessa, 22 Italiassa).

Virulotyyppitystä eli virulenssigeenien etsimistä genomista harjoiteltiin kurssilla sekä CGE:n VirulenceFinder-työkalun (kuva 1) että ARIES:in Virulotyper-työkalun (kuva 2) avulla. VirulenceFinder perustuu blastn-hakuun CGE:n virulenssigeenitietokannasta. Se hyväksyy syötteenä sekä koontisekvenssejä (engl. contig) fasta-tiedostoina että lukusekvenssejä (engl. read) fastq-tiedostoina. Mikäli syötteenä käyttää lukusekvenssejä, työkalu suorittaa *de novo*-koonnin joka tapauksessa ennen blastn-hakua. Virulotyper-työkalu sen sijaan ottaa syötteenä lukusekvenssejä, jotka se rinnastaa CGE:n tietokannan referenssisekvensseihin (engl. mapping, bowtie2-algoritmi). Rinnastus on nopeampaa kuin blastn-haku. Tulokseksi Virulotyper näyttää parhaan osuman kustakin geenistä (tietokannassa useita alleeleja) ja rinnastuksen kattavuuden (engl. coverage). Mikäli kattavuus on alle 10, geenin läsnäoloon kannattaa suhtautua varauksella. Kattavuus vaihtelee kuitenkin geenikohtaisesti; joidenkin genomialueiden sekvensointi on hankalampaa ja kattavuus jää aina matalaksi (esim. liikkuvat geneettiset elementit kuten *stx*). Lisäksi DNA-eristysmenetelmä vaikuttaa siihen, todetaanko plasmideissa sijaitsevia geenejä. EURL:n kokemuksen mukaan DNA-eristys pylväskitillä säilyttää plasmidit hieman paremmin kuin DNA-eristys saostusmenetelmällä. "Ei todettu" -tuloksiin on siis suhtauduttava varauksella. Virulotyyppityksen varmuutta voi lisätä käyttämällä kumpaakin menetelmää (algoritmia) samalle datalle. Lisäksi hakuarvoja väljentämällä voi hakea esimerkiksi geenejä, joiden sekvensointi on onnistunut vain osittain.

Galaxy CGE Server Center for ... Center for ... Center for ... Center for ... Center for ... Center for ... Center for ... Center for ... Center for ...

https://cge.cbs.dtu.dk/cgi-bin/webface.fcgi?jobid=5762983700002D6AA562063D

VirulenceFinder-1.5 Server - Results

1574578.cgebase.cbs.dtu.dk

SETTINGS:
Selected %ID threshold: 90.00

Virulence - E. coli						
Virulence factor	%Identity	Query/HSP length	Contig	Position in contig	Protein function	Accession number
<i>eae</i>	99.93	2820 / 2820	out_69	10801..13620	Intimin	ECU59503
<i>iss</i>	100.00	294 / 294	out_128	13365..13658	Increased serum survival	CP001665
<i>stx2A</i>	100.00	960 / 960	out_170	1384..2343	Shiga toxin 2, subunit A, variant a	AF525040
<i>cif</i>	99.88	849 / 849	out_158	1508..2355	Type II secreted effector	AY128535
<i>espJ</i>	100.00	654 / 654	out_99	1529..2182	Prophage-encoded type III secretion system effector	AB303060
<i>espA</i>	100.00	579 / 579	out_69	16311..16889	Type III secretions system	FM201463
<i>katP</i>	100.00	2211 / 2211	out_160	1683..3893	Plasmid-encoded catalase peroxidase	AB011549
<i>ehxA</i>	100.00	2997 / 2997	out_104	1746..4742	Enterohaemolysin	HM138194
<i>espB</i>	100.00	945 / 945	out_69	18065..19009	Secreted protein B	AF054421
<i>espF</i>	99.84	624 / 624	out_69	20181..20804	Type III secretion system	AF116900
<i>nleA</i>	100.00	1323 / 1323	out_31	2053..3375	Non-LEE encoded effector A	AM422003
<i>efa1</i>	100.00	9672 / 9672	out_69	21739..31410	EHEC factor for adherence	AJ277443
<i>stx2B</i>	100.00	270 / 270	out_170	2355..2624	Shiga toxin 2, subunit B, variant a	AF525040
<i>nleB</i>	100.00	990 / 990	out_69	34011..35000	Non-LEE encoded effector B	AF453441
<i>lpfA</i>	100.00	573 / 573	out_15	37201..37773	Long polar fimbriae	AP010953
<i>espP</i>	100.00	3903 / 3903	out_145	527..4429	Extracellular serine protease plasmid-encoded	GG259888
<i>nleC</i>	100.00	987 / 987	out_8	69368..70354	Non-LEE encoded effector C	AP010960
<i>toxB</i>	100.00	9501 / 9501	out_110	7480..16980	Toxin B	AB456530
<i>tir</i>	99.88	1617 / 1617	out_69	8521..10137	Translocated intimin receptor protein	AB428060

stx - Holotoxins						
Virulence factor	%Identity	Query/HSP length	Contig	Position in contig	Protein function	Accession number
<i>stx2</i>	100.00	1241 / 1241	out_170	1384..2624	O157 SF-258-98, variant a	AF524844

extended output

Results as text Results tab separated Hit in genome sequences Virulence gene sequences

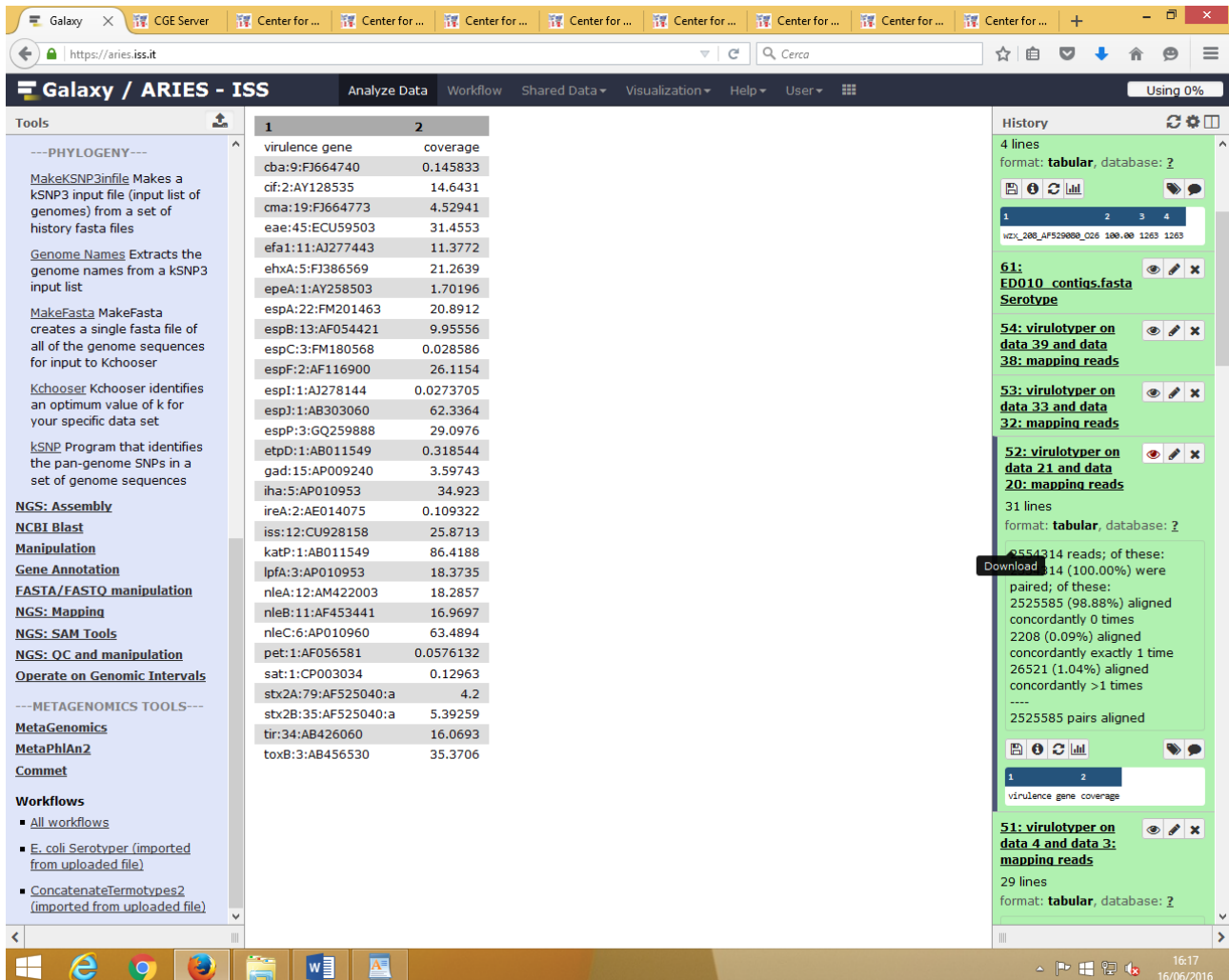
Selected %ID threshold: 90.00 %

Selected minimum length: 60 %

Input Files: ED180_contigs.fasta

Windows taskbar: 16:14 16/06/2016

Kuva 1. VirulenceFinder (CGE) -tulokset kannasta ED180.



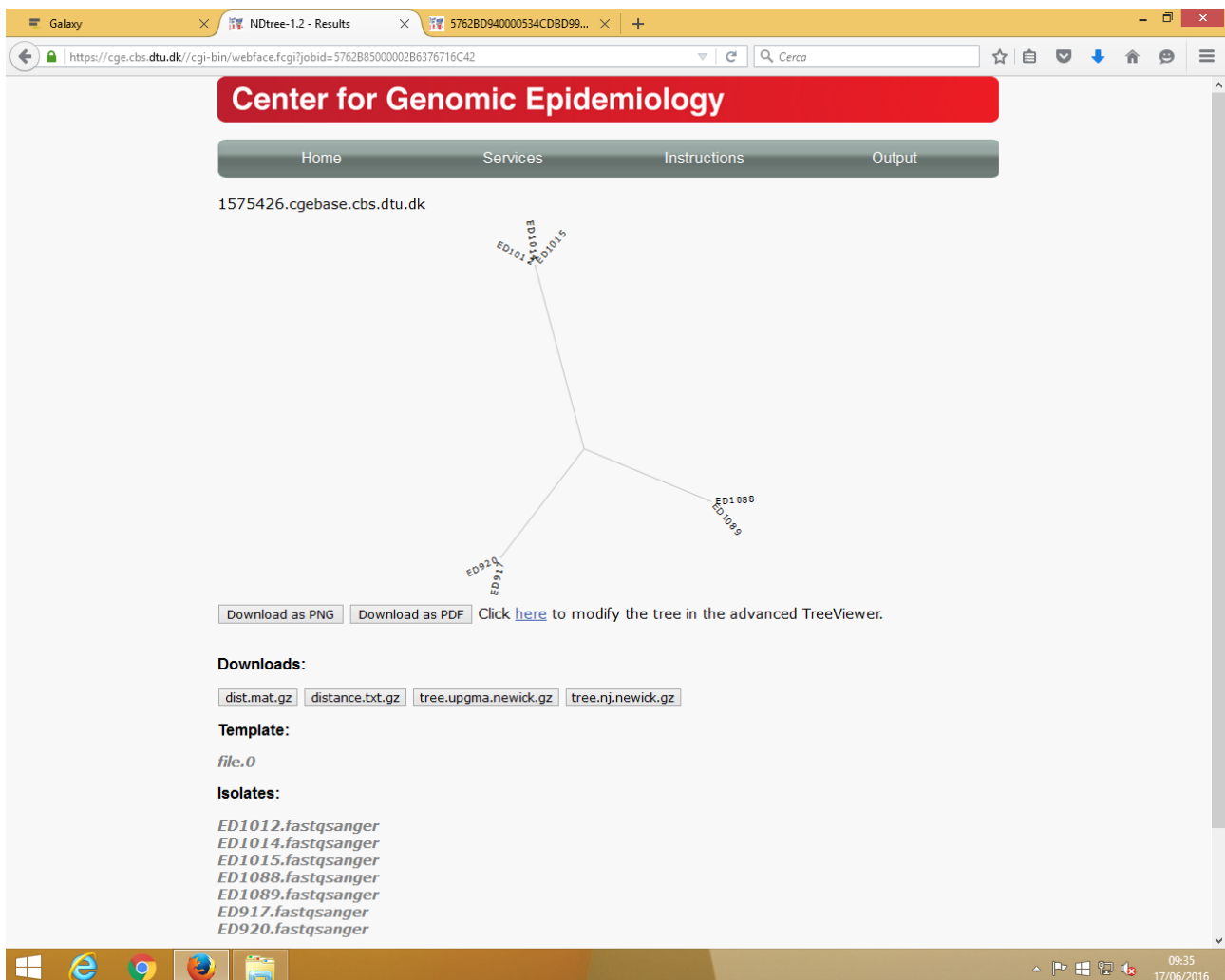
1	2
virulence gene	coverage
cba:9:FJ664740	0.145833
cfi:2:AY128535	14.6431
cma:19:FJ664773	4.52941
eae:45:ECU59503	31.4553
efa1:11:AJ277443	11.3772
ehxA:5:FJ386569	21.2639
epeA:1:AY258503	1.70196
espA:22:FM201463	20.8912
espB:13:AF054421	9.95556
espC:3:FM180568	0.028586
espF:2:AF116900	26.1154
espI:1:AJ278144	0.0273705
espJ:1:AB303060	62.3364
espP:3:GQ259888	29.0976
etpD:1:AB011549	0.318544
gad:15:AP009240	3.59743
iha:5:AP010953	34.923
ireA:2:AE014075	0.109322
iss:12:CU928158	25.8713
katP:1:AB011549	86.4188
lpfA:3:AP010953	18.3735
nleA:12:AM422003	18.2857
nleB:11:AF453441	16.9697
nleC:6:AP010960	63.4894
pet:1:AF056581	0.0576132
sat:1:CP003034	0.12963
stx2A:79:AF525040:a	4.2
stx2B:35:AF525040:a	5.39259
tir:34:AB426060	16.0693
toxB:3:AB456530	35.3706

Kuva 2. Virulotyper (ARIES) -tulokset kannasta ED180.

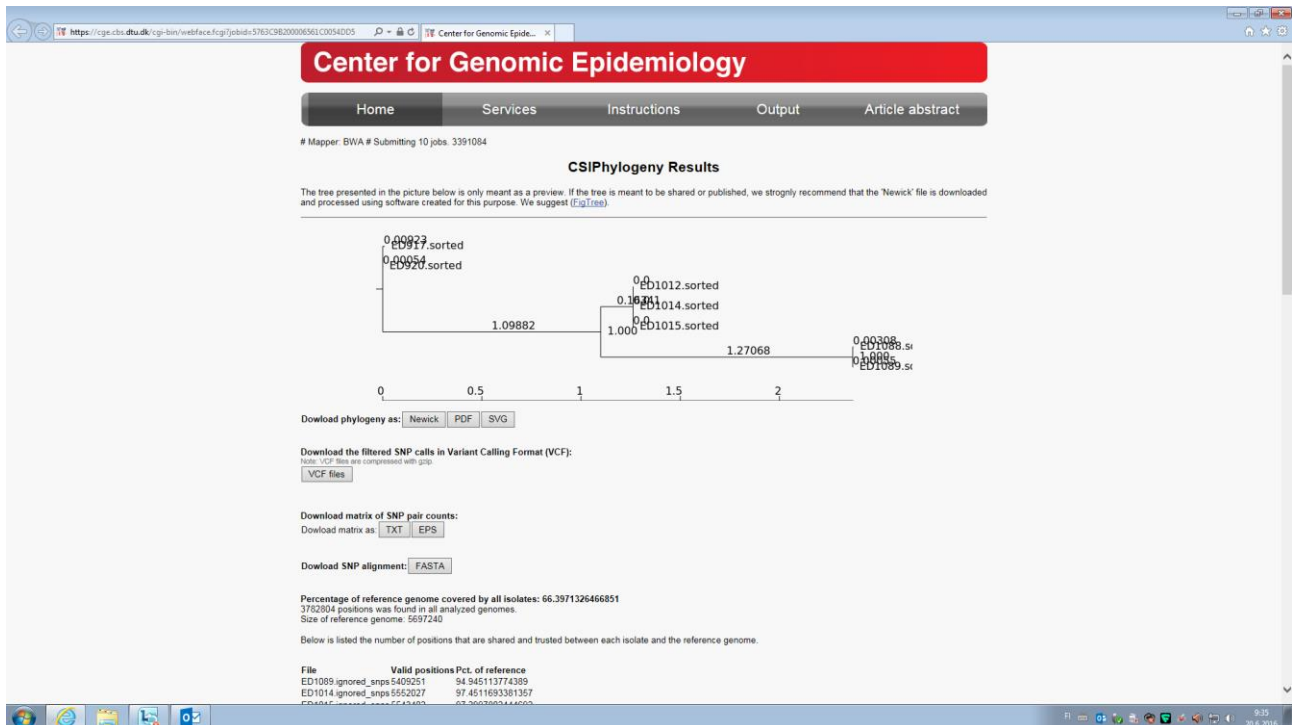
Serotyyppitykseen käytettiin CGE:n SerotypeFinder-työkalua ja ARIES:in Serotyper-työkalua. Kummatkin työkalut perustuvat blastn-hakuun samasta CGE:n tietokannasta ja käyttävät syötteenä fasta-koontisekvenssejä. SerotypeFinderiin voi lisäksi syöttää lukusekvenssejä, mutta työkalu suorittaa niille ensin *de novo*-koonnin.

wgSNP-analyysia eli bakteerigenomien fylogeneettistä analyysiä SNP:ien perusteella kokeiltiin kahden CGE:n työkalun (NDtree, CSI Phylogeny) ja ARIES:in ksnp3-työkalun avulla. CGE:n snpTree-työkalu ei ollut käytettävissä kurssipäivänä. CGE:n työkalut perustuvat referenssigenomiin, johon tutkittavien kantojen lukusekvenssit rinnastetaan. Näin ollen referenssin valinta vaikuttaa oleellisesti tulokseen. CGE:n työkaluista yksinkertaisin ja samalla suljetuin systeemi on NDtree (kuva 3), jossa käyttäjä valitsee itse ainoastaan referenssigenomin, eikä voi vaikuttaa muihin parametreihin. SnpTree sallii käyttäjän valita vähimmäiskattavuuden ja etäisyyden SNP:en välillä, mutta ei tee laadunvarmistusta. CSI Phylogeny (kuva 4) sen sijaan tekee myös laadunvarmistusta, mikäli syötteenä on lukusekvenssejä (hyväksyy myös koonteja). Lisäksi se käyttää fylogenian määrittämiseen vain SNP-sijainteja, joiden laatu kantajoukossa on riittävän hyvä. Se hyväksyy syötteenä myös kantajoukkoja, joista osa on sekvensoitu eri sekvensaattorilla. Kaikki CGE:n työkalut laskevat ja vi-

sualisoivat UPGMA-puun. Neighbor joining -puu on myös saatavilla, mutta sitä varten tulostiedosto täytyy ladata ulkoiseen visualisointiohjelmaan (esim. FigTree). EURL:n kokeilussa NDtree osoittautui epäherkäksi, eikä onnistunut erottelemaan epidemiaan kuulumattomia STEC-kantoja epidemiakannoista. Puusta ei saanut selville etäisyyksiä SNP:ien lukumäärinä. CSI Phylogeny erotteli epidemiakannat muista, muttei onnistunut erottelemaan kantoja saman epidemian sisällä.

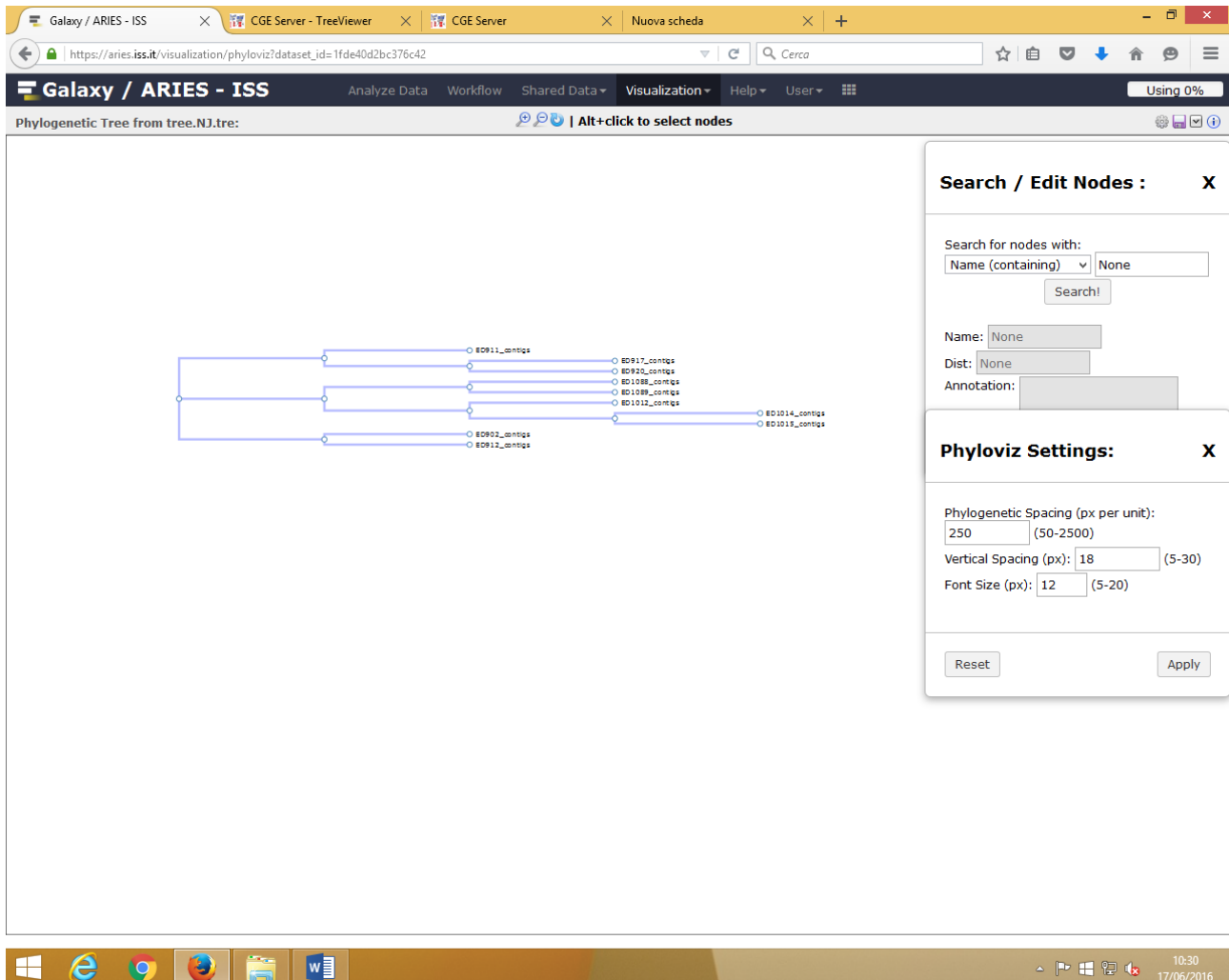


Kuva 3. NDtree-tulos.



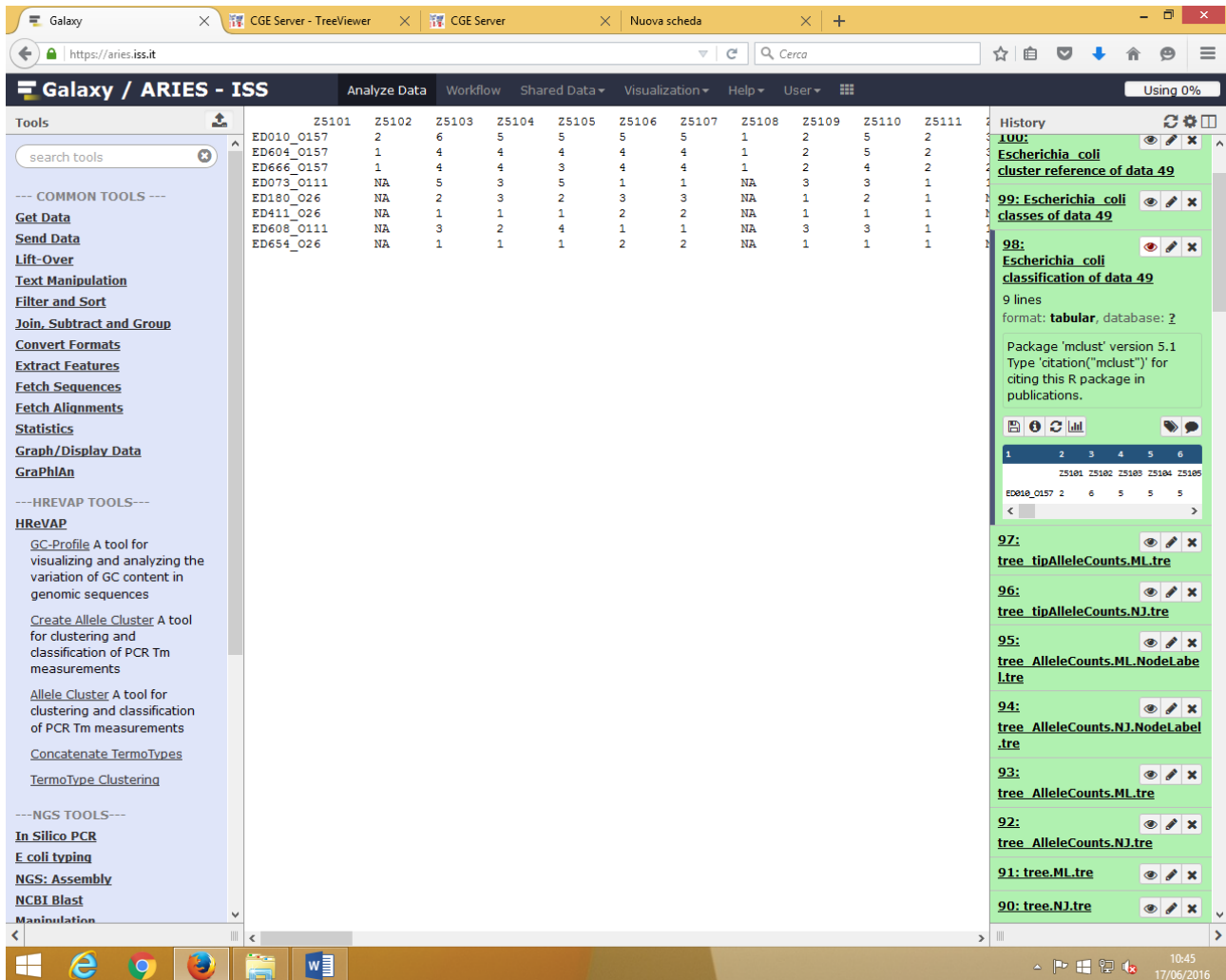
Kuva 4. CSI Phylogeny -tulos. Tuloksista (tulossivulla alempana) näkee lisäksi hylättyjen SNP-sijaintien lukumäärät.

ARIESin ksnp3 (kuva 5) perustuu referenssigenomin sijaan k-mereihin eli työkalu etsii fasta-koontisekvensseistä k:n pituisia nukleotidisekvenssejä, joista se sitten etsii SNP:ejä. Näin työkalu muodostaa SNP-pangenomin tutkittavalle kantajoukolle ja etsii fylogeneettisiä puita (neighbor joining ja maximum likelihood) näiden SNP:ien perusteella. Ennen ksnp-työkalun käyttöä luodaan syötiedostot MakeKSNP3infile- ja MakeFasta-työkalujen avulla sekä optimoidaan k (eli maksimoidaan uniikkien k-merien osuus keskipituudessa koontisekvenssissä) Kschooser-työkalun avulla. Tulospuun voi visualisoida suoraan ARIES:issa Phyloviz-työkalun avulla tai ladata ulkoiseen visualisointiohjelmaan. Ilman referenssigenomia toimivasta menetelmästä on etua erityisesti STEC-genomien fylogeniaa arvioitaessa, sillä STEC-kantojen monimuotoisuus hankaloittaa referenssigenomin valintaa. Varsinkaan täydellisiä, suljettuja genomisekvenssejä ei ole saatavilla harvinaisemmille serotyypeille. SNP:eihin ja referenssigenomiin perustuvat menetelmät soveltuvat parhaiten keskenään samankaltaisten kantojen fylogeneettiseen analyysiin.



Kuva 5. Ksnp3 (ARIES) -tulos, Phyloviz-visualisointi.

HReVAP on EURL:n kehittämä, ARIES:ssa saatavilla oleva analyysimenetelmä virulenssigeenien (91 geeniä, jotka sijaitsevat LEE, OI-57 ja OI-122 -saarekkeissa) alleelien määrittämiseen ja STEC-kantojen fylogeneettiseen ryhmittelyyn alleeliprofiilin perusteella (<http://www.ncbi.nlm.nih.gov/pubmed/26941726>). Menetelmä perustuu alleelien erilaisiin PCR-sulamislämpötiloihin ja on kaksivaiheinen: ensin luodaan alleelitaulukko (kuva 6) sulamislämpötilojen perusteella (työkalu: Allele Cluster) ja sitten suoritetaan alleeliprofiilien klusterointi (työkalu: TermoType Clustering, kuva 7). Alleelitaulukon alleeliluokat voidaan joko luoda tutkittavalle kantajoukolle alusta asti (työkalu: Create Allele Cluster), käyttää aiemmin luotua luokittelua tai systeemireferenssiä. Klusterointi tuottaa neighbor joining -puun. EURL:n kokeiluissa menetelmä on onnistuneesti erotellut seroryhmät toisistaan ja erottelukyky on riittänyt myös kantojen välisten erojen havaitsemiseen seroryhmien sisällä. Erottelukyky on ollut ksnp3-analyysin tasoa. Tosin HReVAP:in erottelukyky on heikompi LEE (*eae*) -negatiivisilla kannoilla.



The screenshot shows the Galaxy web interface with the HReVAP tool output. The main panel displays a table with the following data:

	Z5101	Z5102	Z5103	Z5104	Z5105	Z5106	Z5107	Z5108	Z5109	Z5110	Z5111
ED010_0157	2	6	5	5	5	5	1	2	5	2	
ED604_0157	1	4	4	4	4	4	1	2	5	2	
ED666_0157	1	4	4	3	4	4	1	2	4	2	
ED073_0111	NA	5	3	5	1	1	NA	3	3	1	
ED180_026	NA	2	3	2	3	3	NA	1	2	1	
ED411_026	NA	1	1	1	2	2	NA	1	1	1	
ED608_0111	NA	3	2	4	1	1	NA	3	3	1	
ED654_026	NA	1	1	1	2	2	NA	1	1	1	

The History panel on the right shows a series of tree visualizations for different data points (90-100), including labels like 'Escherichia coli cluster reference of data 49' and 'classification of data 49'. A small preview table is visible for data point 98:

	Z5101	Z5102	Z5103	Z5104	Z5105
ED010_0157	2	6	5	5	5

Kuva 6. HReVAP-alleelitalulukko kahdeksalle kannalle (rivi: kanta, sarake: geeni, solu: alleelinumero). Esimerkiksi geenillä Z5102 on kyseisessä kantajoukossa 6 alleeliluokkaa, alleelinumerot 1–6.

