

ExMEGA, 효과적인 Metagenome 분석 솔루션

1. 시퀀싱 기반 Metagenomics 의 개요

미생물 연구는 오랜 시간 동안 과학의 주요 분야 중 하나로 주목받아왔지만 우리가 알고 있는 미생물의 다양성은 전체의 일부에 불과했다. 미생물의 약 99%는 기존에 알려진 방법으로 배양하기 어려울 뿐만 아니라 배양 독립적인 기술로는 미생물 균집과 그 구성원에 대한 세부 정보를 파악하는 데 한계가 있었기 때문이다. 그러나 Next-Generation Sequencing(NGS) 기술의 도입과 whole-genome shotgun (WGS) sequencing 기술의 등장으로 다양한 환경에 포함된 미생물 균집의 유전체(Metagenome) 구조를 빠르게 분석할 수 있게 되면서 기존의 한계를 극복함과 동시에 시퀀싱 기반 연구에 혁명적인 변화를 일으켰다.[1] 특히, 미국 국립보건원의 Human Microbiome Project (HMP)와 Metagenomics of the Human Intestinal Tract(MetaHIT) 같은 성공적인 국제 마이크로바이옴 프로젝트들이 주목할 만한 성과를 거두면서 NGS 기반의 미생물 연구(Metagenomics)를 황금기로 이끌어주는 계기가 되었다. [2]



그림 1. 국제 마이크로바이옴 프로젝트 [3]

그러나, 기술의 급격한 발전에 따른 대량의 데이터와 NGS 기기의 특성으로 인한 다양한 오류는 metagenomics 분석에서 데이터 품질 관리의 중요성을 더욱 부각시켰다. 이 중 오류 수정(error correction)과 같은 방법이 특히 중요하게 여겨지면서 생물정보학 도구의 중요성이 점점 커지게 되었다.[1] Metagenome 데이터는 다양한 미생물들의 유전 정보를 포함하고 있기 때문에, 그 안에서 특정 미생물의 유전 정보를 분류하거나 분석하는 것은 굉장히 복잡하다. 최근에는 이러한 복잡한 과정을 지원하기 위한 다양한 생물정보학 도구들이 개발되었으며, Raw data 전처리부터 계통분류, 유전체 재구성 등 다양한 단계에서 활용되고 있다. 이러한 고도화된 도구들의 도입은 극한 환경의 미생물 균집을 연구하거나 메타게놈을 유전자나 대사 경로에 연결하는 등의 복잡한 연구를 가능하게 만들고 있다. 이런 복잡한 과정을 거치면서도, 메타게놈 데이터의 효율적 처리와 다양한 분석을 위해 계속해서 새로운 생물정보학 도구들이 개발되고 있다는 것은 주목할 만하다.

1-1. Shotgun Metagenome sequencing and 16s rRNA sequencing

Shotgun Metagenome sequencing 방법은 전체 genome 정보를 획득할 수 있기 때문에 종(species)정보 뿐만 아니라 strain 정보까지 분류학적 정보를 획득할 수 있다. 하나의 샘플에서 다양한 미생물, plasmid, phage, virus, 곰팡이와 같이 다양한 종류의 분석이 가능하며, 유전자의 유무나 기능에 대한 분석도 가능하다. 그러나 host DNA 를 제거하기 어려워 불필요한 sequencing 비용이 발생할 수 있다. 즉 Shotgun Metagenome sequencing 은 많은 분류학적 정보를 얻을 수 있으나, 분석하는데 비용과 시간이 많이 소요된다.

반면 16s rRNA sequencing 의 경우 분류학적으로 의미 있는 16s rRNA 영역을 선택적으로 증폭하여 분석한다. 미생물의 16s rRNA 유전자 영역만을 선택적으로 분석하기 때문에 host DNA 의 오염을 피할 수 있다. 16s rRNA sequencing 은 16s rRNA 영역의 일부만을 분석하기 때문에 분류학적으로는 속(genus) 단계까지 분석 가능하지만, Shotgun Metagenome sequencing 보다 빠르고 합리적인 비용으로 분석이 가능한 장점을 가지고 있다.

2. ExMEGA (Excel based MetaGenome Analysis) Tools

㈜이바이오젠에서 개발 및 보유하고 있는 ExMEGA 는 엑셀 기반으로 실행되며, 다양한 환경내 존재하는 미생물 군집의 생태학적 의미를 분석하는 도구로 활용되고 있다. 또한, Metagenome 데이터에 대한 다양한 분석 기능과 데이터 시각화를 용이하게 해주는 프로그램으로, 데이터 분석과 엑셀 사용에 익숙하지 못한 연구자들도 쉽게 활용할 수 있다. ㈜이바이오젠은 연구자의 요구사항을 지속적으로 반영하여 프로그램을 업데이트하고 있다. 본문에서는 ExMEGA 의 분석 기능과 용어, 기타 분석에 대해 간략히 설명하고자 한다.

ID	ASV	Kingdom	Phylum	Class	Order	Family	Genus	Species	A / B	A / C	B / C
1	00086a129b3dc9e8641815d44e8791d9	d_Bacteri	p_Firmicut	c_Clostridi	o_Clostridi	f_Clostridi	g_Clostridi		0.558031	-0.765062	0.207
2	003986eaa110161172a7ec41b5d5c817	d_Bacteri							-0.041443	-0.078945	0.120
3	00408e6774d91032a45f0b773795c810	d_Bacteri	p_Firmicut	c_Clostridi	o_Lachno	f_Lachnos	g_UCS-1-2		0.23321	-0.440241	0.207
4	004323e46d3b6749f80276f5fe7379bf	d_Bacteri	p_Proteo	c_Alphap	o_Rhodos	f_uncultur	g_uncultur		-0.041443	0.058382	-0.016
5	005c1dbc469b1ba50c1e9d4dc34557fa	d_Bacteri	p_Bacter	c_Bacter	o_Bacter				-0.041443	0.058382	-0.016
6	005e0000f4d659ab67404084dd4af50b	d_Bacteri	p_Firmicut	c_Clostridi	o_Lachno	f_Lachnos	g_Eisenb		-0.526017	-2.294205	2.820
7	006858cc0d2251544b863550d9d067d0	d_Bacteri	p_Bacter	c_Bacter	o_Bacter				-0.041443	0.058382	-0.016
8	008bb2345da6bfa88ab4245db0d74ba	d_Bacteri	p_Firmicut	c_Clostridi	o_Lachno	f_Lachnos	g_Lachnos		-0.297408	-0.165588	0.462
9	00acb08edd1212bdcaef50ecccfd86	d_Bacteri	p_Firmicut	c_Clostridi	o_Lachno	f_Lachnos			-0.865489	-0.165588	1.031
10	00ce8cc6ad0083f7a7cef371b71f8e83	d_Bacteri	p_Firmicut	c_Bacilli	o_Erysipel	f_Erysipel	g_Faecalit		-0.39643	-0.165588	0.562
11	010ddb7bc606ff747373ce07e940ef6a	d_Bacteri	p_Bacter	c_Bacter	o_Bacter				-0.041443	-0.028262	0.069
12	01a26370b5ca9ef7fa78e10095224af8	d_Bacteri	p_Firmicut	c_Clostridi	o_Oscillo	f_Oscillos	g_UCG-00		-0.041443	0.578567	-0.537
13	020f4a9ac9ec34524f8b9a9ec5d6c89b	d_Bacteri	p_Bacter	c_Bacter	o_Bacter	f_Muribac	g_Muribac		-0.338506	-0.165588	0.504
14	021f18a677c19ad0726e7732c25118d2	d_Bacteri	p_Firmicut	c_Clostridi	o_Lachno	f_Lachnos			-0.041443	0.058382	-0.016
15	0223944a7c43dc735783210b2da5f5be	d_Bacteri	p_Firmicut	c_Clostridi	o_Oscillo	f_Ruminoc			0.445035	-0.652066	0.207
16	024cb0459c40803f8cb8db2132ea0e3f	d_Bacteri	p_Firmicut	c_Clostridi	o_Lachno	f_Lachnos	g_Blautia		-0.877277	-0.65488	1.532
17	024ff8f3223f4aacb020e035325f1262	d_Bacteri	p_Firmicut	c_Clostridi	o_Oscillo	f_Eubact	g_Eubact		-0.041443	0.687594	-0.646
18	0277ac1c33495192c73ebd5f8f2dd07c	d_Bacteri	p_Bacter	c_Bacter	o_Bacter	f_Muribac			-0.041443	-0.078945	0.120
19	0277d12d8f01b31d1d2545a6131ef91e	d_Bacteri	p_Bacter	c_Bacter	o_Bacter	f_Muribac	g_Muribac		0.074731	-4.210253	4.135
20	029003bbdd79ec46617325b33949443d	d_Bacteri	p_Firmicut	c_Bacilli	o_Lactob	f_Enteroc			0.478418	-0.685449	0.207

그림 2. 엑셀 기반 ExMEGA 분석 프로그램

2-1. 엑셀내 용어 설명

ASV(Amplicon Sequence Variant)	유사한 DNA 서열 그룹 단위. 즉, Taxon을 구성하는 최소단위
OTU(Operational Taxonomic Unit)	운영 분류 단위: 97% 유사성을 갖는 DNA 서열 그룹
Taxonomy	각 ASV에 대한 7가지 Taxon 정보
CLR(Central Log Ratio)	ANCOM 결과로 나온 중심로그비율
W	ANCOM결과에 따른 통계적 유의성을 나타내는 수치이며, 값이 높을수록 유의미함.
SigW (Significant W)	W결과에 따른 유의미한 수준에 따라 TRUE/FALSE로 반환
Relative Abundance	Feature Count 로 계산된 샘플/그룹 별 비율.
Absolute abundance	생태계의 단위 부피에서 관찰할 수 없는 실제 분류군의 풍부함.
Feature Count	각 ASV에 매핑 된 서열 수
BlastAnnotation	각 ASV의 시퀀스에 대한 Blast 정보

ASV 는 dada2 알고리즘으로 생성되며, 기존 OTU 방법보다 뛰어난 해상도를 보여준다. 더불어, 최근에는 Relative abundance 를 통한 분석만으로는 Metagenome 생태계를 비교하기 어려운 문제로 다양한 통계방식(ANCOM 등)을 이용하여 Absolute abundance 를 추정하고 있다.[4]

2-2. Taxonomic classification

ExMEGA 는 사용자가 원하는 Taxon(분류군)에 따라 필터링하여 데이터를 확인할 수 있는 기능을 제공한다. 이 기능을 통해 Phylum, Class, Order, Family, Genus, Species 와 같은 다양한 분류 수준을 이해하고 분석할 수 있다. 뿐만 아니라, 필터링된 데이터를 Bar plot 과 pie chart 형태로 시각화 하여 정보를 더 명확하게 전달할 수 있도록 지원한다. Metagenome 의 Clustering 최소 단위인 ASV(Amplicon Sequence Variant)는 정해진 분류군을 직접 적으로 볼 수 없기 때문에 실제 Taxon 별로 합산한 상대적 비율을 확인하는 것이 중요하다.

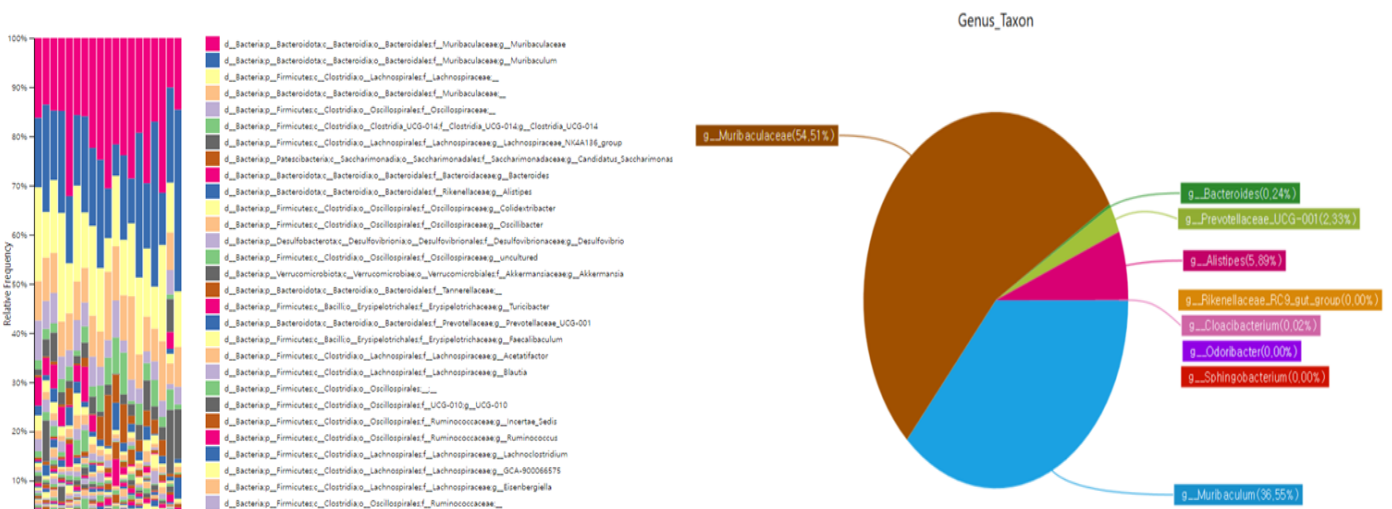


그림 3. ExMEGA Bar plot 과 Pie chart [4]

2-3. DA (Differential abundance) analysis

DA analysis는 주요 통계분석 중 하나로 어떤 특정 변수와 함께 변동하는 미생물 분류군의 풍부도를 찾아내는 것을 목표로 한다. 이렇게 찾아낸 미생물 분류군은 질병 메커니즘에 대한 생물학적 통찰력을 제공할 수 있으며, 질병 예방, 진단, 치료를 위한 바이오 마커로 활용할 수 있다. [5] 해당 분석 또한 ExMEGA를 통해 분석이 가능하며, Volcano plot, Venn Diagram 등 다양한 시각화 이미지를 제공한다.

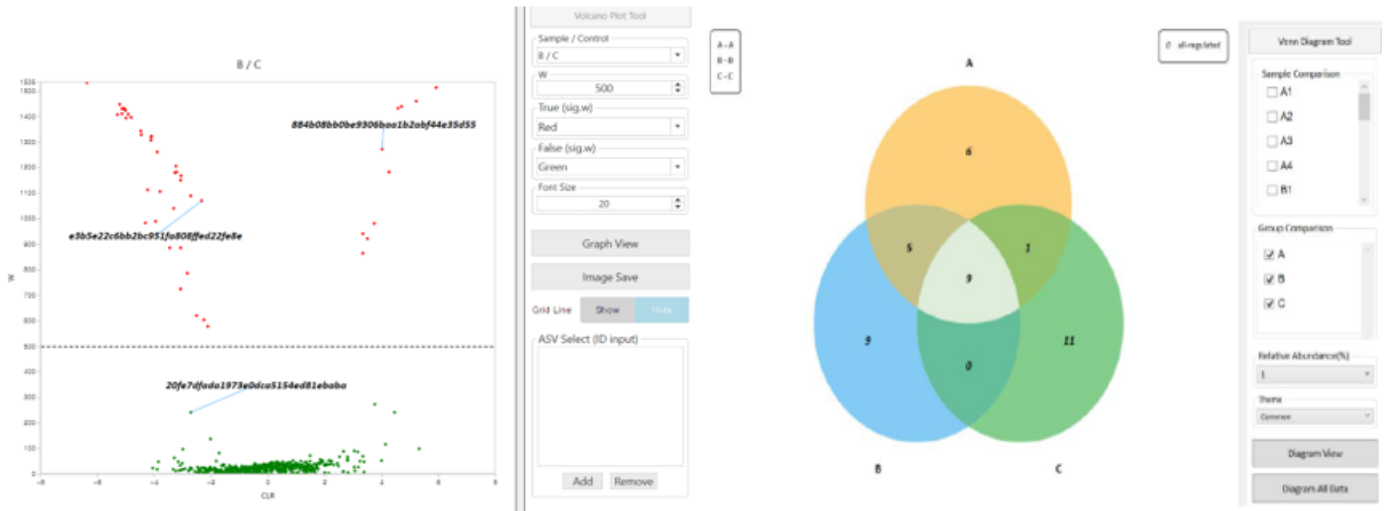


그림 4. ExMEGA Volcano plot, Venn Diagram [4]

2-4. Significant Taxon

아래 스크린샷과 같이 ExMEGA에서는 그룹 별로 비교한 결과에서 유의하게 발현 차이가 나는 분류군을 필터링 할 수 있는 Significant Taxon 도구를 사용할 수 있다. 연구자는 Relative Abundance(상대 풍부도)와 W 값(통계적 유의성)을 기준으로 필터링하여 각 샘플 및 그룹에 대한 조건을 설정하여 분석 결과를 확인할 수 있다. W 값의 경우 값이 높을수록 해당 분류군의 발현차이가 통계적으로 유의미하다는 것을 나타낸다.

Filter: 145		Taxonomy							CLR			
ID	ASV	Kingdom	Phylum	Class	Order	Family	Genus	Specie	A / B	A / C	B / C	
21	19	d__Bacteri	p__Bactero	c__Bactero	o__Bactero	f__Muribac	g__Muribac	s__uncultur	0.074731	-4.210253	4.135	
57	55	08c81d1e3e56673edbd8dbd851a76314	d__Bacteri	p__Bactero	c__Bactero	o__Bactero	f__Muribac	.	0.70749	-0.456576	-0.250	
165	163	1e4fd308329bf0c0adda47c475b6f449	d__Bacteri	p__Firmicut	c__Clostridi	o__Oscillo	f__Rumino	g__Rumino	s__uncultur	2.29258	-1.212805	-1.079
173	171	1f517eb21ccde60de20d3896cfe289f	d__Bacteri	p__Firmicut	c__Clostridi	o__Lachno	f__Lachno	s__.	-0.646465	-0.983086	1.629	
213	211	258520298a5dd98fcc0c46aa0e3c02f7	d__Bacteri	p__Bactero	c__Bactero	o__Bactero	f__Muribac	g__Muribac	s__uncultur	-1.231085	-1.923833	3.154
222	220	26ce56c14635ece701475d32241f18d6	d__Bacteri	p__Bactero	c__Bactero	o__Bactero	f__Muribac	g__Muribac	s__Muribac	0.007134	-0.526613	0.519
237	235	28e99fc648df0be4567b433d614d0752	d__Bacteri	p__Bactero	c__Bactero	o__Bactero	f__Muribac	g__Muribac	s__Muribac	-0.056787	-0.47028	0.527
262	260	2df6c5f78cc097c773d7304ffa17e110	d__Bacteri	p__Bactero	c__Bactero	o__Bactero	f__Muribac	.	0.703621	-0.586351	-0.11	
284	282	3297352312c9b8893a0f9d14eed34b1c	d__Bacteri	p__Bactero	c__Bactero	o__Bactero	f__Muribac	.	-0.166427	-4.297688	4.464	
306	304	35455b695cb7d46fee33151a36fa1b78	d__Bacteri	p__Bactero	c__Bactero	o__Bactero	f__Muribac	g__Muribac	s__uncultur	-1.162694	0.28152	0.881
363	361	3d1dd855a1ef476cf4376494988b66f0	d__Bacteri	p__Bactero	c__Bactero	o__Bactero	f__Muribac	.	1.823967	0.543646	-2.367	
408	406	43ad07aa6643c9dd6d7fd103956c27d0	d__Bacteri	p__Bactero	c__Bactero	o__Bactero	f__Muribac	g__Muribac	s__unidenti	4.917619	0.10165	-5.019
430	428	4785099617e43380e570b85986573b60	d__Bacteri	p__Proteob	c__Alphapr	o__Rhizobii	f__Xanthob	g__Bradyrh	.	0.23321	-0.266955	0.033
464	462	4e499198c57e2ee2063149e9929f1c48	d__Bacteri	p__Bactero	c__Bactero	o__Bactero	f__Muribac	g__Muribac	s__uncultur	2.159526	-0.210797	-1.948
471	469	5003076c835bd08848de4c613cb95cfd	d__Bacteri	p__Bactero	c__Bactero	o__Bactero	f__Muribac	g__Muribac	s__uncultur	-1.998394	1.483133	0.515
472	470	5018f85d7c7c85a225228a19324c9b1	d__Bacteri	p__Bactero	c__Bactero	o__Bactero	f__Muribac	.	-0.26984	-3.487325	3.757	
490	488	530bac5584ac375e0648d99df3b9ccd9	d__Bacteri	p__Bactero	c__Bactero	o__Bactero	f__Muribac	.	-1.458617	-1.884461	3.343	
494	492	5350bf519978ad7b0228cfc91c856d6f	d__Bacteri	p__Firmicut	c__Clostridi	o__Oscillo	f__Rumino	s__.	0.23321	-0.078945	-0.154	
497	495	53b4abdc8f2dd4701f6d319cfb3ca43b	d__Bacteri	p__Bactero	c__Bactero	o__Bactero	f__Muribac	.	2.743704	2.295577	-5.039	
514	512	577a13f0d0d73525dbc12fe21a58b80	d__Bacteri	p__Firmicut	c__Clostridi	o__Clostrid	f__Clostrid	g__Candida	s__Candida	-1.503855	-1.287554	2.791

DA Analysis

Significant Taxon

Relative Abundance(%)
0.00

W
1

Sample / Control

A

A / B

A / C

B / C

All

Pie Chart

Volcano Plot

그림 5. ExMEGA Significant Taxon screenshot

3. ExMEGA GraphicPlus

위에서 설명한 시각화 이미지 외에도 ExMEGA GraphicPlus 를 통한 추가적인 분석이 가능하다. 대표적으로 Krona Chart, PCoA (Principal Coordinates Analysis), Clustering Heatmap 으로 표현 가능하다. 모든 그래프는 단일 샘플 비교(Single Data)와 그룹 비교 (Group Data) 가 모두 가능하며 W 값을 적용하여 분석된다. 만약, 통계분석을 원하지 않을 경우에는 Relative Abundance 만으로도 분석 가능하다.

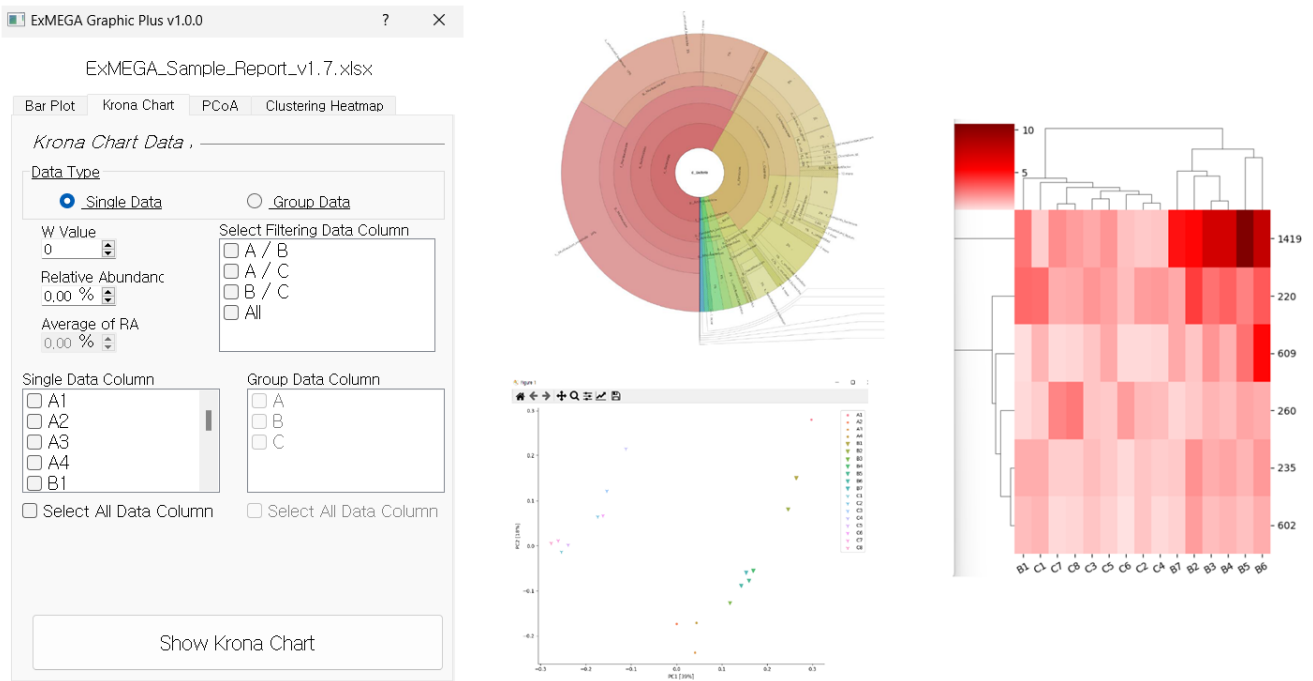


그림 6. ExMEGA GraphicPlus - Krona Chart, PCoA, Clustering Heatmap[4]

4. E-biogen’s Metagenome Service

이바이오젠에서는 Bacteria 의 V3-V4 region, Fungi 의 ITS region 을 증폭하여 실험하는 16s rRNA amplicon sequencing 을 통해 토양, 공기, 물, 조직, 분변시료 등의 다양한 환경에서 채취되는 샘플 내에 어떠한 미생물 (Bacteria, Fungi 등) 들이 얼마만큼 존재하는지 분석할 수 있는 Metagenome & Microbiome 서비스를 제공하고 있다. 해당 서비스를 진행할 경우 위에서 언급한 ExMEGA&GraphicPlus 와 함께 분석 데이터를 제공한다.

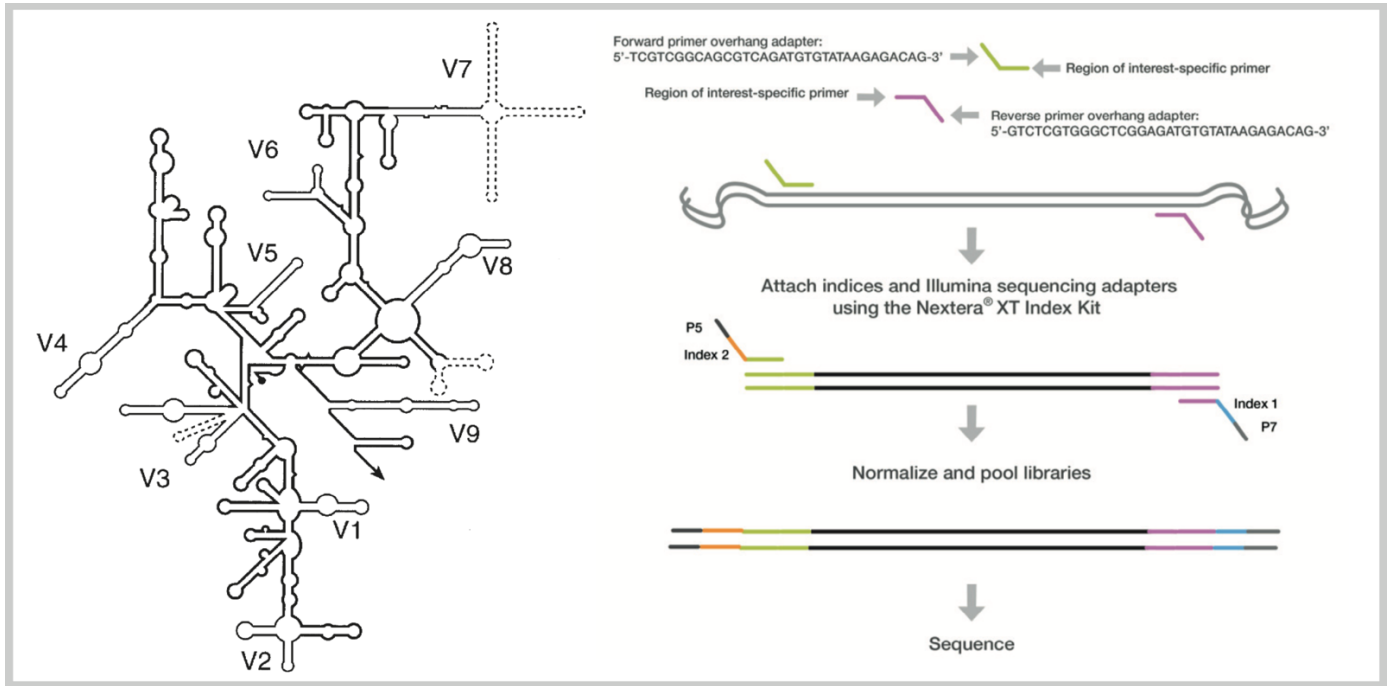
4-1 Service info

Sample requirement	>0.1ng gDNA
Library method	16S rRNA Amplicon Library for illumina
NGS run format	MiSeq, PE300
Data yield	~50,000 read/sample
Turnaround time	Library QC 통과 후 4 주 이내
Sample type	gDNA, fecal, soil, water, feedstuff, tissue, saliva 등

4-2 Metagenome Service Process



4-3 Library Construction Workflow



<참고문헌>

1. Jünemann, S., Kleinbölting, N., Jaenicke, S., et al (2017). Bioinformatics for NGS-based metagenomics and the application to biogas research. *Journal of Biotechnology*, 261, 10–23.
2. Yong-Xin Liu., Yuan Qin., Yang Bai., et al (2021). A practical guide to amplicon and metagenomic analysis of microbiome data. *Protein Cell*, 315–330.
3. 인체 마이크로바이옴 연구개발 동향, KHIDI 전문가 리포트, 4.
4. EBIOGEN, ExMEGA v.10 & Data Analysis User Manual
5. Lu Yang., Jun Chen. (2022). A comprehensive evaluation of microbial differential abundance analysis methods: current status and potential solutions, *Microbiome* 10, 130