

# Use of Near-Neighbor PCR to Close Scaffold Gaps in Microbial Genomes

A. Christine Munk, Yan Xu, Avinash Kewalramani, Riley Arnaudville, Roxanne Tapia, Thomas S. Brettin, C. S. Han  
Los Alamos National Laboratory, Los Alamos, N. M.

## ABSTRACT #FF006

One of the challenges of finishing microbial genomes is large numbers of uncloned regions (scaffold gaps). Gaps with no clone links require many expensive PCR reactions to connect scaffolds. Ordering and orienting scaffolds with no clone links is difficult unless a closely-related genome is available to identify possible links. This abstract describes near-neighbor PCR (nnPCR), a method of closing these 'scaffold' gaps with a minimum of PCR reactions using closely related finished genomes.

In the manual process of near-neighbor PCR, scaffold ends are identified visually using the Consed program's Assembly view and map of scaffolds. Blast is performed using as query a fasta file containing alternatively 1) the sequence of all unordered contigs <2kb and >10 reads, or 2) the final 10kb of each scaffold. A closely-related-genome downloaded from Genbank is used as subject. The blast output is parsed and a tab-delimited file is produced which can be opened as a spreadsheet. Links between contigs are identified and PCR primers are chosen and paired according to the observed links. PCR is performed and products are end-sequenced and assembled to close scaffold gaps. If PCR products are >1500 bp, they are shattered, subcloned and subclones are end-sequenced and assembled to make a consensus sequence to close scaffold gaps.

In one microbial finishing project with 14 uncloned gaps considered, 9 out of 14 PCR reactions were successful, and 9 gaps were closed. In another project with 23 uncloned gaps, 8 out of 19 PCR reactions were successful and 5 gaps were closed. This can be compared to combinatorial PCR, which would require 378 and 703 PCR reactions respectively.

Software to automate and improve this procedure has been developed and is currently being implemented in production. Improvements are currently being tested. This nnPCR software uses the Consed autoreport function to generate a file which contains a map of contigs in scaffolds, and creates a fasta file for the blast query with sequence from the contigs at scaffold ends. nnPCR uses megablast to search all finished genomes currently in Genbank for hits within the same organism within a 20 kb range. The software chooses primers and pairs them, and submits a work order directly to the lab.

For the project with 14 uncloned gaps considered, 12 out of 14 PCR reactions chosen by the nnPCR software were successful. For the project with 23 uncloned gaps, 3 of 6 PCR reactions were successful.

Several microbial genomes have been significantly improved using near-neighbor PCR. This method has significantly reduced the number of PCR reactions that would be required if combinatorial PCR were necessary. Countless hours of finishers' time have been spent on more productive efforts.

US Department of Energy's Office of Science, Biological and Environmental Research Program and Los Alamos National Laboratory under contract No. W-7405-ENG-36

## Acknowledgement:

We would like to thank the Intelligence Technology Innovation Center for funding to sequence the genomes used in this analysis.

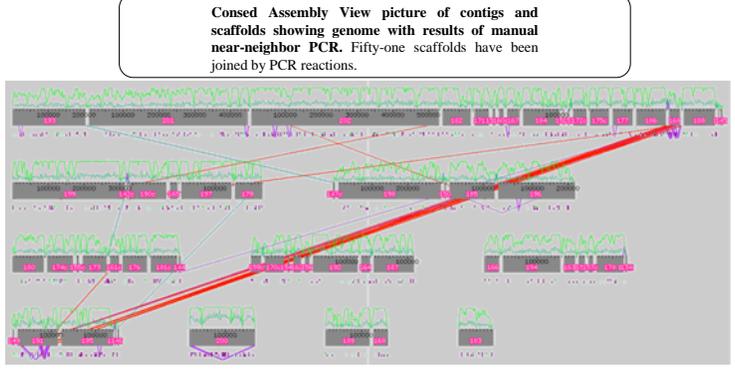
Consed Assembly View picture of contigs and scaffolds showing draft sequence plus 2 rounds of automatic finishing of a *Bacillus anthracis* genome. Each gray bar is a contig. Intercontig lines represent clones with end sequences in different contigs. Adjacent contigs with clones connecting them make a scaffold. Scaffolds with no clone links are not ordered or oriented.



## PCR Results for nnPCR and Manual Process

Paired contigs	manual product	manual product	nnPCR product	nnPCR product	both products	none good	failed
B102-4493	45MC	1400	WZL,WZV	1600;1500	x	x	
B102-4493	45MC	3200	XAL,XAC	3800;3800	x	x	
B147-4071	45ME	700	XAT,XAU	1000;900	x	x	
B171-4018	45MF	1500	none chosen	none chosen	x	x	
B110-4098	45MG	1500	none chosen	none chosen	x	x	
B110-4098	45MH	800	fail	fail	x	x	
B110-4098	45MI	1500	WZD	2100	x	x	
B140-4161	45ML	2800	XAB,XAC	3400	x	x	
B147-4084	45MJ	750	fail	fail	f	x	
B182-4016	45MK	1100	fail	fail	f	x	
B118-4143	45ML	2700	XAJ,XAK	failed	f	x	
B151-4123	45ML	1700	XAJ,XAK	failed	f	x	
B112-4172	45MN	400	XBI	3000	x	x	
B117-4171	45MO	600	XBI	3000	x	x	
B147-4084	45MP	800	fail	fail	x	x	same as above
B148-4181	45MQ	1350	WZB,WZC	1200;1500	x	x	
B182-4016	45MR	1100	WZL	1200	x	x	
B118-4143	45MS	1600	WZV,WZG	400;1800	x	x	
B152-4016	45MT	1100	WZM	1700	x	x	
B14-4101	45MU	1000	XAL,XAM	800;900	x	x	
B159-4088	45MV	800	fail	fail	x	x	
B159-4088	45MW	f	fail	fail	x	x	
B189-4163	45MX	1200	yes	f	x	x	
B189-4163	45MY	1300	fail	fail	f	x	
B189-4163	45MZ	600	WZU,WZV	800;800	x	x	
B189-4163	45NA	1800	WZU,WZV	800;800	x	x	
B189-4163	45NB	1600	WZJ	1600	x	x	
B179-4018	45NC	800	fail	fail	x	x	
B114-4141	45ND	6300	none chosen	none chosen	x	x	
B114-4141	45NE	6300	none chosen	none chosen	x	x	
B119-4180	45NF	6200	none chosen	1100;1100	x	x	
B147-4144	45NG	f	fail	fail	x	x	
B189-4163	45NH	1200	yes	f	x	x	
B148-4181	45NI	1400	none chosen	6000;5000	x	x	
B189-4163	45NJ	3800	WZC,WZD	5000;5000	x	x	
B102-4493	45NK	550	yes	yes	x	x	
B189-4163	45NL	1200	yes	yes	x	x	
B147-4088	45NM	2800	yes	yes	x	x	
B189-4163	45NO	3000	yes	yes	x	x	
B189-4163	45NP	3000	yes	yes	x	x	
B189-4163	45NQ	f	fail	fail	x	x	
B148-4181	45NR	3000	yes	yes	x	x	
B114-4141	45NS	1100	XBF,XBG	900;900	x	x	
B153-4016	45NT	1300	WZL,WZM	1500;1500	x	x	
B158-4157	45NU	1250	none chosen	none chosen	x	x	
B157-4018	45NV	1200	none chosen	none chosen	x	x	
B178-4118	45NW	550	none chosen	none chosen	x	x	
B110-4104	45NX	2300	none chosen	none chosen	x	x	
B118-4158	45NY	1300	fail	fail	x	x	
B189-4163	45NZ	2000	XAF	750	x	x	
B118-4158	45OA	2500	WZL,WZM	4500;4000	x	x	
B153-4016	45OB	1000	XBD,XBE	900;900	x	x	
B173-4018	45OC	7000	yes	yes	x	x	
B141-4071	45OD	1250	yes	yes	x	x	
B188-4018	45OE	7000	yes	yes	x	x	
B118-4158	45OF	4500	none chosen	none chosen	x	x	
B188-4171	45OG	9000	WYX,WYY	2800;2800	x	x	
B117-4118	45OH	4500	none chosen	none chosen	x	x	
B189-4163	45OI	4000	yes	yes	x	x	
B189-4163	45OJ	2000	none chosen	none chosen	x	x	
B189-4163	45OK	2000	none chosen	none chosen	x	x	
TOTALS		36	15	8			

Manual near-neighbor PCR blastParser results. Command-line Blast was performed using a query file containing the sequence of all unordered contigs less than 2kb in length with greater than 10 reads. The finished genome sequence for *Bacillus anthracis* Ames strain was downloaded from Genbank and formatted to use as subject file. Blast output was parsed to output the tab-delimited file shown on the left, ordered by query contigs. The file was copied and re-sorted to represent finished genome order. These files are manually searched to identify links between scaffolds. PCR primers are chosen based on these searches using Consed's primer-picking program. A request for PCR reactions is typed and submitted to the informatics system and forwarded to the lab team.



## Observations and Conclusions

The more closely related microbes are phylogenetically, the more help near neighbors are in connecting sequence scaffolds.

Although the above figures show that additional scaffold links may be found using the manual process, nnPCR software has been recently optimized to pick primer pairs for more scaffold gaps.

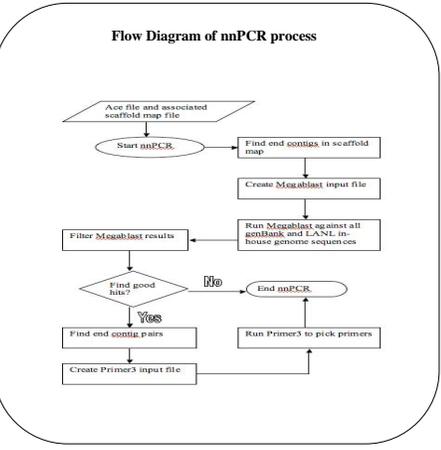
Use of near-neighbor microbial genomes to order and orient scaffolds can significantly reduce the number of PCR reactions necessary to close a genome.

The nnPCR program also reduces the amount of the human finisher's time required to close scaffold gaps by automating the process of identifying scaffold links and submitting work lists.

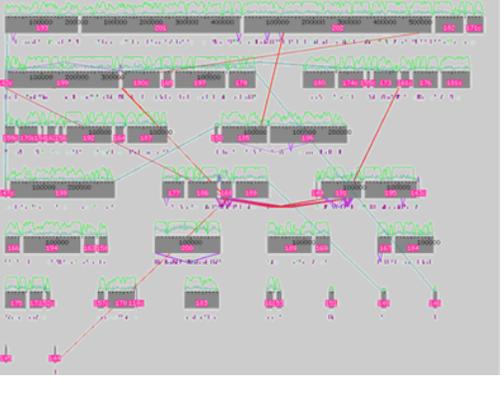
NnPCR can be used with good results as early in the finishing process as after automatic duplication resolution.

## nnPCR Results for Various Microbes

PROJECT	# un-captured gaps	#pcr chosen?	# of gaps closed by nPCR
Yersinia pestis YPC	17	10	6
Bacillus megaterium BMA	61	46	36
Bacillus anthracis BAA	46	28	23
Clostridium botulinum CLJ	26	6	4
Clostridium botulinum CLM	25	6	4
Escherichia coli ECD	38	23	21
Burkholderia pseudomallei BLJ	20	18	17
Vibrio furnessii VFA	52	0	0
Clostridium botulinum CLG	33	0	0



Consed Assembly View picture of contigs and scaffolds showing same genome with results of nnPCR software. Thirty-six scaffolds have been joined by PCR reactions.



NnPCR was performed at various points early in the finishing process for several projects. The number of primer pairs chosen when nnPCR is run immediately following 2 rounds of dupFinisher<sup>®</sup> does not vary significantly from the number of pairs chosen after 1 or 2 cycles of autoFinish, or autoFinish plus some manual finishing. \*automated duplication resolution

PROJECT	Min # scaffolds	nnPCR pairs chosen	Max # scaffolds	nnPCR pairs chosen
Bacillus anthracis BAA	50	46	63	47
Bacillus megaterium BMA	55	46	66	42
Yersinia pestis YPA	17	11	27	18
Yersinia pestis YPK	11	8	15	10
Francisella tularensis FTB	8	6	11	11
Vibrio cholerae VCA	37	19	49	22
Burkholderia pseudomallei BLJ	31	18	36	15