

MDPI

Article

Comparative Analysis and Ancestral Sequence Reconstruction of Bacterial Sortase Family Proteins Generates Functional Ancestral Mutants with Different Sequence Specificities

Jordan D. Valgardson ^{1,†}, Sarah A. Struyvenberg ^{1,‡}, Zachary R. Sailer ^{2,3,§}, Isabel M. Piper ^{1,||}, Justin E. Svendsen ^{1,2}, D. Alex Johnson ^{1,¶}, Brandon A. Vogel ¹, John M. Antos ¹, Michael J. Harms ^{2,3} and Jeanine F. Amacher ^{1,*}, □

- Department of Chemistry, Western Washington University, Bellingham, WA 98225, USA; jvalgard@stanford.edu (J.D.V.); sastruyvenberg@msn.com (S.A.S.); isabel_piper@berkeley.edu (I.M.P.); jsvends2@uoregon.edu (J.E.S.); dajohnso@caltech.edu (D.A.J.); vogelb2@wwu.edu (B.A.V.); antosj@wwu.edu (J.M.A.)
- Department of Chemistry and Biochemistry, University of Oregon, Eugene, OR 97403, USA; zachsailer@gmail.com (Z.R.S.); harms@uoregon.edu (M.J.H.)
- Institute of Molecular Biology, University of Oregon, Eugene, OR 97403, USA
- * Correspondence: amachej@wwu.edu; Tel.: +1-360-650-4397
- † Current address: Department of Chemical and Systems Biology, Stanford University, Palo Alto, CA 94305, USA.
- ‡ Current address: Lumen Biosciences, Seattle, WA 98103, USA.
- § Current address: Apple Inc., Cupertino, CA 95014, USA.
- Current address: Department of Chemistry, University of California, Berkeley, CA 94720, USA.
- ¶ Current address: Department of Bioengineering, CalTech, Pasadena, CA 91125, USA.

Abstract: Gram-positive bacteria are some of the earliest known life forms, diverging from gramnegative bacteria 2 billion years ago. These organisms utilize sortase enzymes to attach proteins to their peptidoglycan cell wall, a structural feature that distinguishes the two types of bacteria. The transpeptidase activity of sortases make them an important tool in protein engineering applications, e.g., in sortase-mediated ligations or sortagging. However, due to relatively low catalytic efficiency, there are ongoing efforts to create better sortase variants for these uses. Here, we use bioinformatics tools, principal component analysis and ancestral sequence reconstruction, in combination with protein biochemistry, to analyze natural sequence variation in these enzymes. Principal component analysis on the sortase superfamily distinguishes previously described classes and identifies regions of relatively high sequence variation in structurally-conserved loops within each sortase family, including those near the active site. Using ancestral sequence reconstruction, we determined sequences of ancestral Staphylococcus and Streptococcus Class A sortase proteins. Enzyme assays revealed that the ancestral Streptococcus enzyme is relatively active and shares similar sequence variation with other Class A Streptococcus sortases. Taken together, we highlight how natural sequence variation can be utilized to investigate this important protein family, arguing that these and similar techniques may be used to discover or design sortases with increased catalytic efficiency and/or selectivity for sortase-mediated ligation experiments.

Keywords: sortases; enzymes; protein engineering; principal component analysis; network analysis; bioinformatics; ancestral sequence reconstruction; evolution



Citation: Valgardson, J.D.;
Struyvenberg, S.A.; Sailer, Z.R.; Piper,
I.M.; Svendsen, J.E.; Johnson, D.A.;
Vogel, B.A.; Antos, J.M.; Harms, M.J.;
Amacher, J.F. Comparative Analysis
and Ancestral Sequence
Reconstruction of Bacterial Sortase
Family Proteins Generates Functional
Ancestral Mutants with Different
Sequence Specificities. *Bacteria* 2022,
1, 121–135. https://doi.org/10.3390/
bacteria1020011

Academic Editor: Bart C. Weimer

Received: 29 April 2022 Accepted: 7 June 2022 Published: 9 June 2022

Publisher's Note: MDPI stays neutral with regard to jurisdictional claims in published maps and institutional affiliations.



Copyright: © 2022 by the authors. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (https://creativecommons.org/licenses/by/4.0/).

1. Introduction

Gram-positive bacteria accounted for 76% of all bloodstream infections in 2000, up from 62% in 1995 [1]. Although varied by region and over time, these numbers have stayed relatively consistent for the past 20 years [2–4]. These organisms are defined in part by their thick peptidoglycan layer as compared to gram-negative bacteria, which

they diverged from roughly 2 billion years ago [1,5,6]. Sortase enzymes are critical for the ability of gram-positive bacteria to attach proteins to the cell exterior, as well as to build the pili [7–10]. Due to this activity, sortases are a potential therapeutic target for antibiotic development, and they are actively-used tools for protein engineering [11,12]. Several of the infections mentioned above are caused by pathogenic *Staphylococci* and *Streptococci*, e.g., *Staphylococcus aureus* and *epidermidis*, and *Streptococcus pneumoniae*, *pyogenes*, and *agalactiae* [1]. Therefore, a greater understanding of proteins from these organisms may prove valuable in the fight against gram-positive bacterial infection.

There are six main classes of sortase (class A-F); the first-characterized and best-studied bacterial sortase is the Class A sortase from *Staphylococcus aureus* (saSrtA) [13]. This enzyme recognizes the Cell Wall Sorting Signal (CWSS) sequence LPXTG, where X = any amino acid. Following cleavage of the initial protein target, an acyl-enzyme intermediate is formed. A secondary substrate then acts as a nucleophile, and a final ligation product is generated [9]. Peptidase activity occurs between the Thr and Gly residues, and positions are defined as P4 = Leu, P3 = Pro, P2 = X, P1 = Thr, and P1' = Gly. Other Class A sortases, e.g., *Streptococcus pyogenes* SrtA (spySrtA), are predicted to contain a closely related recognition mechanism, and our group recently showed that recognition of the P1' residue is partially mediated by residues in the $\beta4-\beta5$ and $\beta7-\beta8$ loops, highlighting the importance of these conserved structural features [14,15].

The catalytic activity of sortases make them an exciting tool in protein engineering, where sortase-mediated ligation (SML) or sortagging techniques are commonly employed to create a variety of products, including the recent development of an in vivo assay using engineered saSrtA to label amyloid-β protein in human cerebrospinal fluid and the implementation of ligation site switching to allow assembly of multiple fragments using a single sortase enzyme, amongst many others [11,16–18]. Despite their uses, sortagging applications are hindered by the poor relative enzymatic efficiency of saSrtA and other naturally occurring sortases studied to date [19–21]. Directed evolution studies performed in 2011 were successful in generating a saSrtA pentamutant (P94R/D160N/D165A/K190E/K196T) with an overall catalytic efficiency increase >100-fold [21]. Engineering of additional variants of saSrtA and other Class A sortases is an area of ongoing work. An example includes the incorporation of two additional mutations to the saSrtA pentamutant at the calcium-binding site, which led to a calcium-independent saSrtA heptamutant [20,22–24]. Other studies use directed evolution or other engineering techniques to alter the substrate specificity of saSrtA, e.g., a recent study that reported an saSrtA variant which recognizes an LMVGG substrate motif in the amyloid- β protein [17].

Variation in substrate selectivity also naturally exists amongst bacterial sortases. Although saSrtA is selective for the LPXTG target sequence, this is not true of all Class A sortases. Work from ourselves and others revealed that other Class A sortases can recognize a variety of amino acids at multiple positions [14,15,25,26]. A complete understanding of the selectivity determinants of these alternate preferences is not known. Furthermore, there are six known classes of sortases (A–F). Many of these classes share a similar recognition motif as Class A sortases, including Classes C-F (Class C: [I/L][P/A]XTG; Class D: LPNTA; Class E: LAXTG; Class F: less is known, but it is likely similar to SrtA, LPXTG) [27,28]. However, the recognition motif of Class B sortases is NP[Q/K]TN [27]. Taken together, we hypothesize that investigating sequence variation of individual classes of sortases, as well as the sortase superfamily, may identify sortases with improved catalytic efficiency and/or unique recognition motifs.

Ancestral sequence reconstruction (ASR) is a powerful technique that combines our growing knowledge of the proteomes of extant organisms with statistical methods in order to predict the sequences of ancestral proteins [29]. These ancestral proteins can then be characterized, providing evolutionary clues to sequence–function relationships in a growing number of protein systems, including classic models, e.g., recent work on the origin of cooperativity in hemoglobin [30]. A number of studies suggest that ancestral proteins are less selective for target ligands and more thermostable than extant sequences [31].

Therefore, we propose that ASR can be used as a method for identifying improved sortase sequences for protein engineering.

Here, we used principal component analysis (PCA) and ASR to study the sortase superfamily and Class A sortase sequence variation, respectively. Using PCA, we show that the main source of natural variation within sortase families occurs in a number of structurally-conserved loops near the active site. Using ASR, we characterized ancestral proteins of the genera *Staphylococus* and *Streptococcus*. While our ancestral *Staphylococcus* protein revealed lower relative activity than saSrtA, the ancestral *Streptococcus* enzyme had the second-highest activity of the four *Streptococcus* SrtA proteins studied in similar experiments [14,15]. Interestingly, the ancestral *Streptococcus* SrtA showed markedly increased activity and P1 promiscuity, as compared to its extant *S. pneumoniae* relative [14,15]. Although ancestral sortases from nodes that included multiple genera were expressed and purified, these enzymes were catalytically inactive, due to a number of potential factors. Overall, our work suggests that the ancestral *Streptococcus* protein was relatively more active as compared to its extant relatives and that the ASR technique provides a viable approach for exploring sequence variation in sortases from the same genera.

2. Results

2.1. Principal Component Analysis (PCA) of Bacterial Sortases

In order to gain a better understanding of global sequence patterns in the sortase superfamily, we used PCA to group and analyze 39,188 sortase sequences from all classes. This work builds off of recent studies that utilized a sequence similarity network to classify sortases [27]. Briefly, we downloaded all sequences annotated as "sortase" from UniProt and aligned them by MAFFT, followed by PCA [32,33]. The amino acids in each sequence were then classified by five parameters: hydrophobicity, disorder propensity, molecular weight, charge, and occupancy (defined as a binary value, where 1 = amino acid and 0 = insertion or deletion (indel) at this position) [34,35]. PCA was then performed on the resulting matrix. For visualization purposes this data was projected onto the first three principal components which describe 42.7% of the total variance (Figure S1a). Additionally, we performed Hierarchical Gaussian Mixture Model clustering of the sortase superfamily, as described in the Materials and Methods. On the entire principal component space we hierarchically fit a two Gaussian mixture model to the data until each subcluster reached a minimum size or the Gaussian mixture modeling process failed to identify two distinct Gaussians [36]. The resulting tree from this process can accurately distinguish the known sortase classes, as well as extract small subclusters of sortases and present them in a readable manner (Figure 1a). We also plotted our PCA using the top three principal components (Figure S1b). For visualization, we ran PCA on a subset of the data, including 9427 sequences that were filtered for low numbers of indels and manually verified (Figure 1b).

This analysis verified previous classifications of sortases based on sequence alignment, network, and phylogenetic tree analyses [27,28,37]. For example, principal component 1 (PC1) separates the sortase F proteins from the rest of the superfamily and PC2 captures the separation between sortase B and the other sortase families, as well as sortase E and sortase A. These analyses allowed us to identify the regions of highest variability within each class based on the parameters defined above. We plotted our data onto previously determined sortase A structures by taking the distance from the centroid for each position in the multiple sequence alignment (Figure S1c). Consistent with expectations, we found that secondary structure elements are highly conserved, including the "sortase fold" β -barrel core and class-specific α -helices (Figures 1c,d and S1d). Additionally, PCA revealed that the highest degree of variability occurs in structurally conserved loops adjacent to the substrate recognition pocket (Figures 1c,d and S1d).

Given that the $\beta6-\beta7$ loop has been shown to be intimately involved in sortase substrate recognition in Staphylococcus aureus SrtA (saSrtA), we were intrigued that PCA revealed similar levels of variability in the $\beta4-\beta5$ and $\beta7-\beta8$ loops [38]. In the case of $\beta7-\beta8$,

the resulting matrix. For visualization purposes this data was projected onto the first principal components which describe 42.7% of the total variance (Figure Additionally, we performed Hierarchical Gaussian Mixture Model clustering of sortase superfamily, as described in the Materials and Methods. On the entire principal component space we hierarchically fit a two Gaussian mixture model to the data each subcluster reached a minimum size or the Gaussian mixture modeling process

to identify two distinct Gaussians [36]. The resulting tree from this process can accu we were also motivated by previously reported mutations in the β7–β8 loop of saSrtA distinguish the known sortase classes, as well as extract small subclusters of sortase that have been shown to dramatically modulate sortase reaction rates [8,21,39]. Indeed, our work confirms that the β7–β8 loop dramatically affects the activity and substrate specificity of a sortase principal components (Figure S1b). For visualization, we ran PCA one subset more promisculate, including 9427 sequences that preferring for proving the process of including 9427 sequences sthat preferring for proving the process of including 9427 sequences sthat preferring for proving the proving figure 1b).

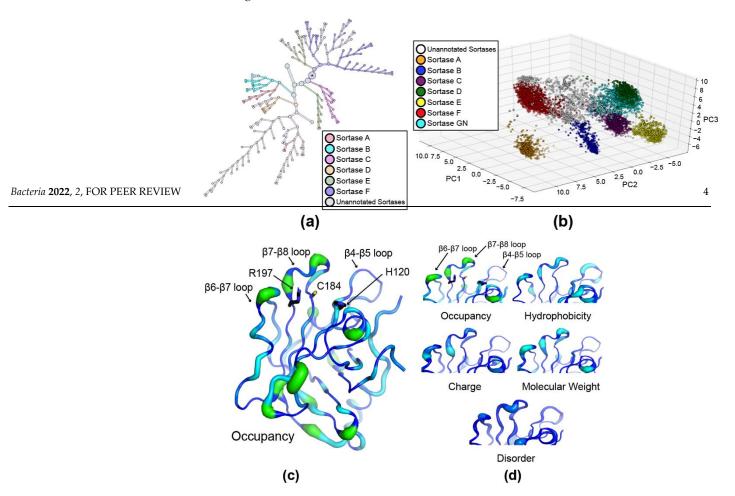


Figure 1. Principal compensations lyie (PCA) of sortage superfamily reveals accurate varied bility stricts transpectations and all different indicting the prof. So the sectast accurately be to a sortage accurately be a sortage accurately accurately be a sortage accurately a

This analysis verified previous classifications of sortases based on sequence alignment, network, and phylogenetic tree analyses [27,28,37]. For example, principal component 1 (Principal phylogenetic tree analyses [27,28,37]. For example, principal component 1 (Principal phylogenetic tree analyses [27,28,37]. For example, principal component 1 (Principal phylogenetic tree analyses [27,28,37]. For example, principal component 1 (Principal phylogenetic tree analyses [27,28,37]. For example, principal component 1 (Principal phylogenetic tree analyses [27,28,37]. For example, principal component 1 (Principal phylogenetic tree analyses [27,28,37]. For example, principal component 1 (Principal phylogenetic tree analyses [27,28,37]. For example, principal component 1 (Principal phylogenetic tree analyses [27,28,37]. For example, principal component 1 (Principal phylogenetic tree analyses [27,28,37]. For example, principal component 1 (Principal phylogenetic tree analyses [27,28,37]. For example, principal component 1 (Principal phylogenetic tree analyses [27,28,37]. For example, principal component 1 (Principal phylogenetic tree analyses [27,28,37]. For example, principal component 1 (Principal phylogenetic tree analyses [27,28,37]. For example, principal component 1 (Principal phylogenetic tree analyses [27,28,37]. For example, principal component 1 (Principal phylogenetic tree analyses [27,28,37]. For example, principal phylogenetic tree analyses [27,28,37]. For example, principal component 1 (Principal phylogenetic tree analyses [27,28,37]. For example, principal component 2 (Principal phylogenetic tree analyses [27,28,37]. For example, principal component 2 (Principal phylogenetic tree analyses [27,28,37]. For example, principal component 2 (Principal phylogenetic tree analyses [27,28,37]. For example, principal component 2 (Principal phylogenetic tree analyses [27,28,37]. For example, principal component 2 (Principal phylogenetic tree analyses [27,28,37]. For example, principal component 2 (Principal phylogenetic

2.2. Ancestral Sequence Reconstruction of Class A Sortases

Building off our PCA analysis of sortase families, we wanted to further explore se-125 quence space in these enzymes by performing ancestral sequence reconstruction (ASR) on Class A sortases. As detailed in the Materials and Methods, ultimately 400 sequences— Artifliting vertes expective sanschib to the chroudbing politing phyting phyting tractic twere—used for all forse-quences used to substitute the sanschib to the chroudbing phyting phyting tractic twere—used for all forse-quences used to substitute the sanschibation of the s

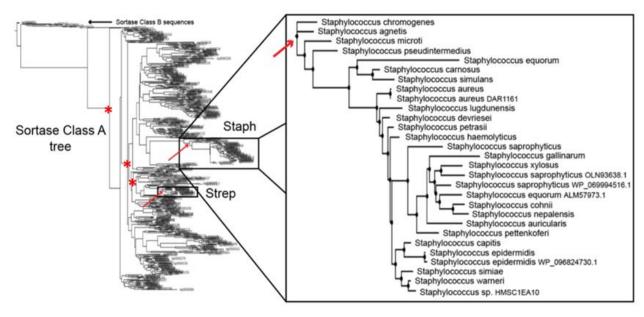


Figure 2. Constructed phylogenetic tree of SrtA sequences used for ancestral sequence reconstruction. The locations of ancistaphsrtA (also, inset) and ancistrepsrtA are indicated by red arrows. More The locations of ancistaphsrtA (also, inset) and ancistrepsrtA are indicated by red arrows. More ancestral proteins are indicated by asterisks; moving from left to right, these are ancinode-408, anancestral proteins are indicated by asterisks; moving from left to right, these are ancinode-408, anancinode-503, and ancinode-547.

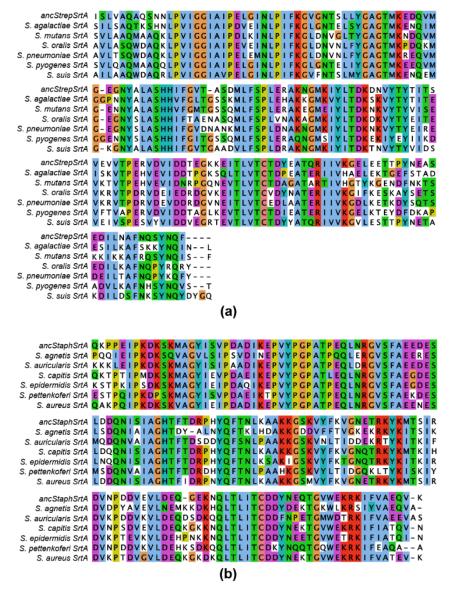
Multiple sequence alignments of our ancestral proteins with representative extant sequences reveals approximate values of our ancestral proteins with representative extant sequences reveals approximate values of 78.3% identifying for one sparkers that with another sequence rand 6,9 fg. identifying an estimptor value to pathers to the other representative extantive extantive and the sequence of the other representative extantive extantive in the sequence of the other representative extantive extantive to an extantive for most of the other representative extantive extantive the aligned erecions for an extantive that aligned erecions for an extantive extantive extantive the aligned erecions for an extantive extantive extantive extantive the aligned erecions for an extantive extant sequences, respectively.

Thabtel1. Sequence identities of aneStaphSrtA and aneStepSFtA with extantequences on interest of identical residues control to be to the training of aneStaphSrtA and aneStepSFtA with extantequences on interesting the sequences of the sequences

(a)	a nesiStap h	s. S_{ig}ngne tis	S. Lu cus eus	S. a b ri eur äeuslari	s S. Espanyitis	SSeppidenmidis	Ssppttenkoferi
a nus Eata ph	χX						
SSagagetestis	68%% (9 7 97/14212)	хх					
s\$afffereus	86%(127/147)	767679(108/142)	χX				
S. auricularis	73% (107/146)	62% (91/146)	70% (101/144) X			
S. auricularis S. capitis	88% (149)146)	62% (91/146) 62% (88/141)	(\$4 ¹ / ₄)2/146	X 68% 68% (96/141)	Х		
S. capitis	88% (129/146)	62% (88/141)	84% (122/146)	(96/141)	X		
S. epidermidis	81% (116/143)	59% (84/142)	78% (114/146)	69% (98/142)	84% (124/147)	Х	
S. pettenkoferi	74% (105/141)	58% (81/140)	67% (96/143)	78% (109/140)	73% (104/143)	71% (101/143)	Х

Bacteria 2022, 1, FOR PEER REVIEW 6 126

S. epidermidis	81% (116/143)	Table 1. <i>Cont.</i> 59% (84/142)	78% (114/146)	69% (98/142)	84% (124/147)	Х	
S. pettenkoferi	74% (105/141)	58% (81/140)	67% (96/143)	78% (109/140)	s ⁷ 3% (104/143)	71% (101/143)	X S. suis
ancStrep SansStretjae	<i>ançStrep</i> (113 X 165)	S. agalactiae X	S. mutans	S. oralis	S. pneumoniae	S. pyogenes	S. suis
S. agalactiae S. mutans S. mutans	68% (513/165) 65% (1069363)	64% (109/169) 64% (109/169)	× _X				
S. oralis S. oralis	69% 69%1(3/13/3)63)	(97/166) (97/166)	65% 6 5 %7(106%)165)	x x			
S. pneumoniae S. pneumoniae	68% 68% (1/1/164) 71%	(95)/167) (95/167)	64% 64% (10%)166) 70%	81% (136/167) 81% (136/167)	X X		
S. pyogenes S. pyogenes	71%4/1/17/4)64)	65% (109/168)	7 0%7(1158) 168)	63% (104/166)	63% (105/166)	×	
<i>§.</i> કૃપાંક	76% 76% <u>x</u> 525(4)64)	558% (288/168)	60% 60%1/401/168)	61%1640(4/046/4065)	622%10303/4606)	63% (104/166)	χХ



Figures3 Multiples opposes and ignores to choose surphy tests in sinit by representatives surphy the process and (M) Streptocous States proteins Multiple be equivered by the proteins of any phyloger constand Streptococus States posterior several a window of the process with this heart in the proteins of any phyloger constand.

To characterize the ancestral proteins, we recombinantly expressed and purified these eTochnestantizethsenvertalspyteistalspyt

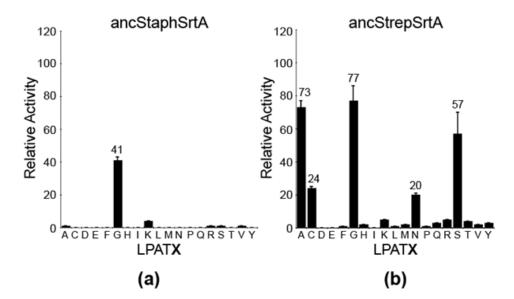


Figure 4. Relative enzyme activities for ancestral proteins. Substrate selectivity data for ancestral proteins and substrate selectivity data for ancestral proteins. Substrate selectivity data for ancestral proteins as a substrate selectivity data for ancestral proteins. Substrate selectivity data for ancestrate and appropriate selectivity data for ancestrate selectiv

2.3.35 & tratedralia A. And glyeseo fo A. Aceses trials & A.A. Porteleis is

Introduct de buttet interpres on blickbanicial datawe a seet de bonology modeling do predict the structure of an acceptably retains 41-143 [ITHe etaplate structure of an acceptably retains 41-143 [ITHE etaplate structure of an acceptably retains 41-143 [ITHE etaplate structure of an acceptably retains a section of the prediction of the interpretable of a structure of a production interior placed in a current of a production of the active acquire at the active acquire active acquire at the active acquire acquire active acquire active acquire acquire

Bacteria **2022**, FOR PEER REVIEW

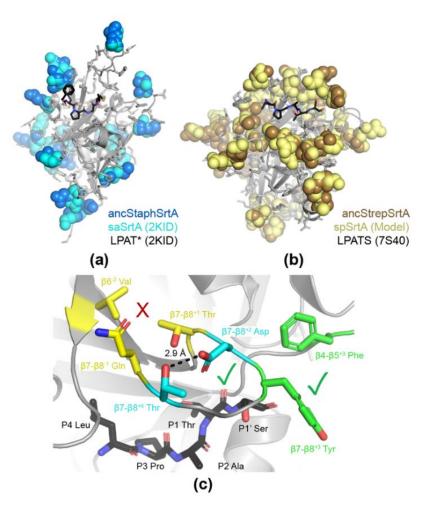


Figure 5.5 Structural comparison of amestral protein models with will diversely to Spread in the tenes will the person of the protein models with the post of the person of the person

Comparison of ancStaphSrtA to saSrtA (PDB ID 2KID) revealed very few changes near the approximate fraction of the first the control of the first fraction of the first fraction

(Figure 4) [38].

In our analysis of ancStrepSrtA and spSrtA, we used a previously generated homology In our analysis of ancStrepSrtA and spSrtA, we used a previously generated homology model of spSrtA, as the only available structures are of a domain-swapped dimer whose ogy model of spSrtA, as the only available structures are of a domain-swapped dimer activity has yet to be confirmed (Figure 5b) [14]. Notably, alignment of our homology whose activity has yet to be confirmed (Figure 5b) [14]. Notably, alignment of our homology model with the predicted structure from the AlphaFold Protein Structure Database reveals or one mean squared deviation (RMSD) for main chain atoms of 0.501 Å (489 atoms), with reveals a root mean squared deviation (RMSD) for main chain atoms of 0.501 Å (489

the largest amount of variation in the $\beta6-\beta7$ and $\beta7-\beta8$ loops (Figure S2) [45,46]. Although these sequences are less similar than the Staphylococcus proteins, we again observe relatively few amino acid substitutions in residues that directly interact with the ligand (Figure 5b). Here, we used the LPATS peptide from spySrtA-LPATS (PDB ID 7S40) for reference. We do observe amino acid variants in the β7–β8 loop residues of ancStrepSrtA as compared to spSrtA that may explain the increased activity of the ancestral protein. As we have previously described, an interaction between the $\beta 6^{-2}$ (or two residues from the C-terminus of the $\beta6$ strand) R184 and two residues in the spSrtA $\beta7$ – $\beta8$ loop, $\beta7$ – $\beta8$ ⁺¹ (or 1 residue C-terminal to the catalytic Cys) E208 and $\beta7-\beta8^{-1}$ (or 1 residue N-terminal to the catalytic Arg) E214, weakens the overall activity of spSrtA [14]. In contrast, spySrtA does not contain this interaction and shows much higher relative activity [15]. AncStrepSrtA contains a $\beta 7 - \beta 8^{+1}$ Thr, $\beta 7 - \beta 8^{-1}$ Gln, and $\beta 6^{-2}$ Val, suggesting that this interaction is also not present in this protein (Figure 5c). We do, however, observe that ancStrepSrtA likely conserves the two favorable interactions previously described that are mediated by $\beta7$ – $\beta8$ loop residues, including an intra-loop hydrogen bond between $\beta 7 - \beta 8^{+2}$ Asp and $\beta 7 - \beta 8^{+6}$ Thr, as well as a hydrophobic interaction between the $\beta7-\beta8^{+3}$ Tyr residue and $\beta4-\beta5^{+3}$ Phe (or three residues C-terminal to the catalytic His) (Figure 5c) [14,15].

2.4. Investigating Ancestral Proteins at Distant Nodes

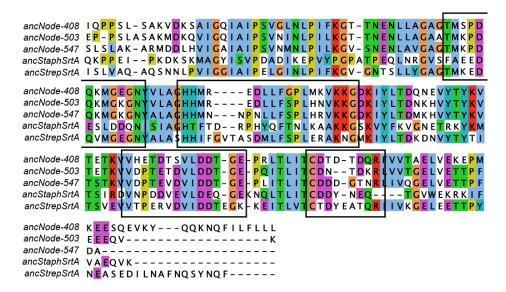
Finally, we wanted to test the activity of ancestral SrtA proteins at more distant nodes in our ASR analyses. We chose three sequences with relatively low sequence identity to ancStaphSrtA and ancStrepSrtA that were also distinct from each other (Figure 2, Table 2). All protein sequences are in the Supplemental Information. We named the proteins for their node characterization in the ASR, ancNode-408, ancNode-503, and ancNode-547.

	ancStaph	ancStrep	ancNode-408	ancNode-503	ancNode-547
ancStaph	Х				
ancStrep	30% (35/117)	Х			
ancNode-408	35% (47/133)	51% (77/151)	Х		
ancNode-503	33% (50/150)	56% (76/136)	78% (156/200)	Х	
ancNode-547	54% (64/118)	59% (85/145)	64% (147/199)	77% (168/199)	Х

Table 2. Comparative pairwise sequence identities of ancestral proteins in this study.

We expressed and purified these proteins as described in the Materials and Methods. Notably, only fractions corresponding to the monomeric peak were retained following size exclusion chromatography, and based on their migration, these proteins are not aggregated and retain a similar radius of gyration as the wild-type proteins (Figure S3). Unfortunately, when evaluated using our FRET-based assay, all three proteins were catalytically inactive for sequences containing P1' Ala, Gly, and Ser residues. Multiple sequence alignment of the ancestral proteins in this study suggests why these proteins may be catalytically inactive (Figure 6). Specifically, the manual refinement of the multiple sequence alignment used for ASR aimed to reduce numbers of gaps in the overall alignment, thereby optimizing alignments in areas of conserved secondary structure elements, e.g., the eight-stranded β -barrel structure conserved in the characterized sortase fold [9,14]. In doing so, we predict this introduced gaps in the structurally-conserved loops near the active site, e.g., the $\beta4-\beta5$ and $\beta7-\beta8$ loops previously mentioned here (Figure 6). The $\beta6-\beta7$ loop appears largely conserved in length, perhaps indicative of a higher degree of length conservation in this structural feature, as well as the $\beta3-\beta4$ loop, which, while spatially more distant from the active site, contains residues previously implicated in ligand recognition (Figure 6) [44].

Bacteria **2022**, 2, FOR PEER REVIEW Bacteria **2022**. 1



The loop lengths of the \$7-\text{B8-loops} of Class A sortases can vary quite dramatically. The loop lengths of the \$7-\text{B8-loops} of Class A sortases can vary quite dramatically. We previously characterized the \$7-\text{B8-loops} of several Class A sortases which varied from 7 residues in \$\frac{\text{Streptococcus}}{\text{sproteins}}\$ proteins to 12 residues in \$\frac{\text{Streptococcus}}{\text{sproteins}}\$ proteins to 12 residues in \$\frac{\text{Streptococcus}}{\text{sproteins}}\$ in \$\frac{\text{sproteins}}{\text{sproteins}}\$ to 12 residues in \$\frac{\text{streptococcus}}{\text{sproteins}}\$ in \$\frac{\text{sproteins}}{\text{sproteins}}\$ in \$\frac{\text{sproteins}}{\t

3. Discussion

In this work, we used bioinformatics tools, IPCA and ASR, to investigate sequence variation in sortases. We used IPCA too dentify soorces of a variation in the sortases performing apitaphiting previous work that the characterized the thief different seastes be dond in this tient to eque properties [127]474 clotically tion found that whithin this ease, dashat gest variationally include and unstructurally colly execute the active site of the east part of the entire active site of the east of the east of the earth of the east of th

These loops were further implicated as sources of relatively high variation in our ASR biochemical studies. Htere, while we were able to express and purify enzymatically active ancestral proteins of the *Staphylococcus* and *Streptococcus* genera, ancestral sequences of nodes that combined multiple genera were catalytically inactive. We predict that this is due to truncations in the 64–655 and 666–657 loops assertesult of manual multiple sequence alignment refriencement iduring ASR REGIBEURO. Desprishes, tois and sugsteen streps from the sequence alignment refriencement iduring ASR REGIBEURO. Desprishes, tois and sugsteen similar similar activity [assats]. Wel als Whater far datest appstraps than a final similar desire of intention of the sequence identity with sags the (109/168 residues), and 63% (105/166) with sags the Aspares of (109/168 residues), and 63% (105/166) with sags the Aspares of (109/168 residues), and 63% (105/166) with sags the Aspares of (109/168 residues), and 63% (105/166) with sags the Aspares of (109/168 residues), and 63% (105/166) with sags the Aspares of (109/168 residues).

While ancestral proteins at deep nodes that included multiple genera as descendants were found to be inactive, the fact that they were able to be expressed and purified using the same methods as those used for extant proteins suggested that the central sortase fold remained intact. Future work to repeat the ASR with careful attention to loop lengths, as well as introduction of extant $\beta4-\beta5$ and $\beta6-\beta7$ loops into ancestral proteins, could provide a means for restoring activity to these enzymes, and may elucidate additional molecular characteristics of the contribution of these individual regions to the activity and selectivity of sortases. Such information would be very useful in future design efforts for sortase enzymes with improved catalytic efficiency or altered specificity. It would also be interesting to perform structural studies on these ancestral proteins, providing insights into potential differences compared to extant proteins with respect to the stereochemistry of target recognition.

There are a number of potential tools that can be used to examine sequence variation in bacterial sortases. Here, we utilized network and evolutionary approaches to investigate natural sequence variation. We argue that with the existence of thousands of sortase enzymes in multiple classes, there is still much to be discovered in extant sortase sequences [27,48]. In addition, directed evolution has proved to be an exciting technique to engineer sortase variation in vitro [17,21,49]. Both approaches, investigating natural sequences, as well as introducing new variation, will allow for a deeper understanding of the sequence determinants of activity and target selectivity, and can profoundly impact the study of the sortase enzyme family, both in protein engineering and for therapeutic uses.

4. Materials and Methods

Principal component analysis (PCA). Initial sequences were obtained from UniProt and an alignment was generated by MAFFT [32,33]. Initially, each sequence was given a score for the number of gaps present for each residue and the filtered alignment was realigned by MAFFT version 7. Subsequent analysis included all sequences without taking gaps into consideration (Figure 2b vs. Figures 2a and S2c). The sortase multiple sequence alignment (MSA) was converted to a tensor of sequences, characterized by MSA position and chemical property of each amino acid. Each amino acid was associated with 4 biochemical traits and a binary trait occupancy, as described. Each trait was normalized to the range from zero to one. In addition, gaps were given the average value of the matrix column with the exception of occupancy so that they would not contribute to variance of the column. Gapped positions were given an occupancy score of zero (for the other chemical properties gapped positions received the average score). After translating the MSA, the resulting tensor was flattened to matrix stacking of the chemical properties and was re-centered so that the matrix had a column-wise mean of zero. Principal component analysis was performed on the matrix by the singular value decomposition algorithm provided in the scikit learn Python package [50]. Clustering was performed by a Gaussian mixture model provided in the scikit-learn 1.1 Python package [50]. Optimal cluster numbers were scored by Bayesian information criterion. Visualization was performed using a script written in Python with matplotlib. Programs were run using default parameters, unless otherwise noted.

Ancestral Sequence Reconstruction. Nonredundant sortase sequences were sourced from the NCBI protein database [51]. Cluster Database at High Intensity with Tolerance program (CD-HIT) was used to filter out highly similar (>95%) identical sequences sourced from NCBI [52,53]. An all-vs-all basic local alignment search tool (BLAST) was used on the remaining sortase sequences, producing a sortase network which informed the assignment of sortase class groups (A–F) by using labeled sortase sequences to assign a class to each grouping [54]. Proteins surrounding the class A group were selected and an additional round of filtering was performed, where all highly similar proteins (>90%) were filtered out via CD-HIT. The remaining pool of sortase sequences was then subjected to alignment by MUltiple Sequence Comparison by Log-Expectation (MUSCLE), and then manually curated to remove any outlying sequences [55]. Seven Class B sortase

sequences (from *Streptococcus suis, Streptococcus oralis, Streptococcus pneumoniae, Staphylococcus aureus, Bacillus anthracis, Listeria monocytogenes*, and *Enterococcus faecalis*) were added to anchor the resulting phylogenetic tree. The final alignment contained a total of 400 sequences. SrtA structures sourced from the PDB database were structurally aligned and sequence similarity between structural sequences (via PDB) and sortase sequences from the multi-sequence alignment (MSA) (via ASR) then informed the true alignment of the MSA. A phylogenetic tree was constructed from the MSA via phyml 3.0 and ancestral sequences were then generated at each node via multi-channel access XML (maxml) [56]. These latter steps were run using a python script. The aLRT values for proteins characterized were ancStaphSrtA = 15.6525, ancStrepSrtA = 13.0091, ancNode-408 = 17.7893, ancNode-503 = 28.8809, and ancNode-547 = 17.5286. Programs were run using default parameters, unless otherwise noted.

Protein expression and purification. Recombinant ancestral proteins (ancStaphSrtA, ancStrepSrtA, ancNode-408, ancNode-503, and ancNode-547) were expressed using Escherichia coli BL21 (DE3) cells in the pET28a(+) vector (Genscript), as previously described [14,15]. Transformed cells were grown at 37 °C in LB media to an OD $_{600}$ 0.6–0.8, followed by induction using 0.15 mM IPTG for 18–20 h at 18 °C. The cells were harvested in lysis buffer [0.05 M Tris pH 7.5, 0.15 M NaCl, 0.5 mM ethylenediaminetetraacetic acid (EDTA)] and whole cell lysate was clarified using centrifugation, followed by filtration of the supernatant. The supernatant was initially purified using a 5 mL HisTrap HP column (Cytiva), with wash [0.05 M Tris pH 7.5, 0.15 M NaCl, 0.02 M imidazole, 0.001 M TCEP] and elution [wash buffer with 0.3 M imidazole] buffers.

Following immobilized metal affinity chromatography, the protein was concentrated using an Amicon Ultra-15 Centrifugal Filter Unit (10,000 NWML) followed by size exclusion chromatography (SEC) using a HiLoad 16/600 Superdex 75 column (Cytiva), with SEC running buffer [0.05 M Tris pH 7.5, 0.15 M NaCl, 0.001 M TCEP]. Purified protein fractions corresponding to the monomeric peak were pooled and concentrated. Purity was assessed using SDS-PAGE. Protein concentrations were determined using theoretical extinction coefficients calculated using ExPASy ProtParam [57]. Protein not immediately used was flash-frozen in SEC running buffer and stored at $-80\,^{\circ}$ C.

Fluorescence Assay for Sortase Activity. Model peptide substrates with the general structure Abz-LPATXG-K(Dnp) (Abz = 2-aminobenzoyl, Dnp = 2,4-dinitrophenyl) were synthesized and purified as previously described [14]. Reactions were analyzed using a Biotek Synergy H1 plate reader as previously described [14,15]. Briefly, reactions were performed a 100 μ L reaction volume consisting of 5 μ M sortase, 50 μ M peptide substrate, 5 mM hydroxylamine nucleophile, and 10% (v/v) 10× sortase reaction buffer (500 mM Tris pH 7.5, 1500 mM NaCl, and 100 mM CaCl₂). The reactions were performed in triplicate and the fluorescence intensity of each well was measured at 2-min time intervals over a 2-h period at room temperature (λ_{ex} = 320 nm, λ_{em} = 420 nm, and detector gain = 75). For each substrate sequence, the background fluorescence of the intact peptide in the absence of enzyme was subtracted from the observed experimental data. Background-corrected fluorescence data was then normalized to the fluorescence intensity of a benchmark reaction between wild-type saSrtA and Abz-LPATGG-K(Dnp), as previously described [14,15].

Structural analyses. Alignments were visualized using AliView [58]. Phylogenetic trees were visualized with FigTree v1.4.3 [59]. Homology modeling was performed using the SwissModel web interface [41,43]. Structural analyses and figure rendering were done using PyMOL. Enzyme assay graphs were prepared using Kaleidagraph. The Streptococcus pneumoniae SrtA structure was downloaded from the AlphaFold Protein Structure Database (entry number Q8DPM3) [45,46].

Supplementary Materials: The following supporting information can be downloaded at: https://www.mdpi.com/article/10.3390/bacteria1020011/s1, Figure S1: Principal Component Analysis (PCA) of sortase superfamily; Figure S2: Structures of spSrtA from the AlphaFold database and homology modeling; Figure S3: Size exclusion chromatography of ancestral SrtA proteins; Recombinant protein sequences used in this study [14,25].

Author Contributions: Conceptualization, J.D.V., Z.R.S., M.J.H. and J.F.A.; methodology, J.D.V., Z.R.S., M.J.H. and J.F.A.; software, J.D.V., Z.R.S. and M.J.H.; validation, S.A.S., I.M.P., J.E.S., D.A.J. and B.A.V.; formal analysis, J.D.V., J.M.A. and J.F.A.; investigation, J.D.V., S.A.S., Z.R.S., I.M.P., J.E.S., D.A.J. and B.A.V.; resources, J.M.A., M.J.H. and J.F.A.; data curation, J.D.V. and J.F.A.; writing—original draft preparation, J.F.A.; writing—review and editing, J.D.V., S.A.S., I.M.P., J.E.S., J.M.A., M.J.H. and J.F.A.; visualization, J.D.V. and J.F.A.; supervision, J.M.A., M.J.H. and J.F.A.; project administration, J.F.A.; funding acquisition, J.M.A., M.J.H. and J.F.A. All authors have read and agreed to the published version of the manuscript.

Funding: J.F.A. and J.M.A. were both funded by Cottrell Scholar Awards from the Research Corporation for Science Advancement. J.F.A. was also funded by NSF CHE-CAREER-2044958. M.J.H. was funded by NSF DEB-CAREER-1844963. In addition, I.M.P. received an Elwha Undergraduate Summer Research Award and D.A.J. received a Joseph & Karen Morse Student Research in Chemistry Fellowship to fund summer research.

Institutional Review Board Statement: Not applicable.

Informed Consent Statement: Not applicable.

Data Availability Statement: Not applicable.

Conflicts of Interest: The authors declare no conflict of interest.

References

- Sizar, O.; Unakal, C.G. Gram Positive Bacteria. In StatPearls; StatPearls Publishing: Treasure Island, FL, USA, 2022.
- 2. Diekema, D.J.; Hsueh, P.-R.; Mendes, R.E.; Pfaller, M.A.; Rolston, K.V.; Sader, H.S.; Jones, R.N. The Microbiology of Bloodstream Infection: 20-Year Trends from the SENTRY Antimicrobial Surveillance Program. *Antimicrob. Agents Chemother.* **2019**, *63*, e00355-19. [CrossRef] [PubMed]
- 3. Maharath, A.; Ahmed, M.S. Bacterial Etiology of Bloodstream Infections and Antimicrobial Resistance Patterns from a Tertiary Care Hospital in Malé, Maldives. *Int. J. Microbiol.* **2021**, 2021, 3088202. [CrossRef] [PubMed]
- 4. Zhu, Q.; Yue, Y.; Zhu, L.; Cui, J.; Zhu, M.; Chen, L.; Yang, Z.; Liang, Z. Epidemiology and microbiology of Gram-positive bloodstream infections in a tertiary-care hospital in Beijing, China: A 6-year retrospective study. *Antimicrob. Resist. Infect. Control* 2018, 7, 107. [CrossRef] [PubMed]
- 5. Vollmer, W.; Blanot, D.; de Pedro, M.A. Peptidoglycan structure and architecture. FEMS Microbiol. Rev. 2008, 32, 149–167. [CrossRef] [PubMed]
- 6. Feng, D.F.; Cho, G.; Doolittle, R.F. Determining divergence times with a protein clock: Update and reevaluation. *Proc. Natl. Acad. Sci. USA* **1997**, *94*, 13028–13033. [CrossRef] [PubMed]
- 7. Marraffini, L.A.; DeDent, A.C.; Schneewind, O. Sortases and the Art of Anchoring Proteins to the Envelopes of Gram-Positive Bacteria. *Microbiol. Mol. Biol. Rev.* **2006**, *70*, 192–221. [CrossRef]
- 8. Ton-That, H.; Mazmanian, S.K.; Alksne, L.; Schneewind, O. Anchoring of surface proteins to the cell wall of Staphylococcus aureus. Cysteine 184 and histidine 120 of sortase form a thiolate-imidazolium ion pair for catalysis. *J. Biol. Chem.* **2002**, 277, 7447–7452. [CrossRef]
- 9. Jacobitz, A.W.; Kattke, M.D.; Wereszczynski, J.; Clubb, R.T. Sortase transpeptidases: Structural biology and catalytic mechanism. *Adv. Protein Chem. Struct. Biol.* **2017**, 109, 223–264.
- 10. Spirig, T.; Weiner, E.M.; Clubb, R.T. Sortase enzymes in Gram-positive bacteria. Mol. Microbiol. 2011, 82, 1044–1059. [CrossRef]
- 11. Antos, J.M.; Truttmann, M.C.; Ploegh, H.L. Recent advances in sortase-catalyzed ligation methodology. *Curr. Opin. Struct. Biol.* **2016**, *38*, 111–118. [CrossRef]
- 12. Zhang, J.; Liu, H.; Zhu, K.; Gong, S.; Dramsi, S.; Wang, Y.-T.; Li, J.; Chen, F.; Zhang, R.; Zhou, L.; et al. Antiinfective therapy with a small molecule inhibitor of Staphylococcus aureus sortase. *Proc. Natl. Acad. Sci. USA* **2014**, *111*, 13517–13522. [CrossRef] [PubMed]
- 13. Ton-That, H.; Liu, G.; Mazmanian, S.K.; Faull, K.F.; Schneewind, O. Purification and characterization of sortase, the transpeptidase that cleaves surface proteins of Staphylococcus aureus at the LPXTG motif. *Proc. Natl. Acad. Sci. USA* 1999, 96, 12424–12429. [CrossRef] [PubMed]
- Piper, I.M.; Struyvenberg, S.A.; Valgardson, J.D.; Johnson, D.A.; Gao, M.; Johnston, K.; Svendsen, J.E.; Kodama, H.M.; Hvorecny, K.L.; Antos, J.M.; et al. Sequence variation in the β7–β8 loop of bacterial class A sortase enzymes alters substrate selectivity. *J. Biol. Chem.* 2021, 297, 100981. [CrossRef]
- 15. Gao, M.; Johnson, D.A.; Piper, I.M.; Kodama, H.M.; Svendsen, J.E.; Tahti, E.; Longshore-Neate, F.; Vogel, B.; Antos, J.M.; Amacher, J.F. Structural and biochemical analyses of selectivity determinants in chimeric Streptococcus Class A sortase enzymes. *Protein Sci.* **2022**, *31*, 701–715. [CrossRef]
- 16. Dai, X.; Böker, A.; Glebe, U. Broadening the scope of sortagging. RSC Adv. 2019, 9, 4700–4721. [CrossRef] [PubMed]

Bacteria **2022**, 1 134

17. Podracky, C.J.; An, C.; DeSousa, A.; Dorr, B.M.; Walsh, D.M.; Liu, D.R. Laboratory evolution of a sortase enzyme that modifies amyloid-β protein. *Nat. Chem. Biol.* **2021**, *17*, 317–325. [CrossRef]

- 18. Bierlmeier, J.; Álvaro-Benito, M.; Scheffler, M.; Sturm, K.; Rehkopf, L.; Freund, C.; Schwarzer, D. Sortase-Mediated Multi-Fragment Assemblies by Ligation Site Switching. *Angew. Chem. Int. Ed.* **2022**, *61*, e202109032. [CrossRef]
- 19. Kruger, R.G.; Otvos, B.; Frankel, B.A.; Bentley, M.; Dostal, P.; McCafferty, D.G. Analysis of the substrate specificity of the Staphylococcus aureus sortase transpeptidase SrtA. *Biochemistry* **2004**, *43*, 1541–1551. [CrossRef]
- 20. Freund, C.; Schwarzer, D. Engineered sortases in peptide and protein chemistry. Chembiochem 2021, 22, 1347–1356. [CrossRef]
- 21. Chen, I.; Dorr, B.M.; Liu, D.R. A general strategy for the evolution of bond-forming enzymes using yeast display. *Proc. Natl. Acad. Sci. USA* **2011**, *108*, 11399–11404. [CrossRef]
- 22. Hirakawa, H.; Ishikawa, S.; Nagamune, T. Design of Ca²⁺-independent Staphylococcus aureus sortase A mutants. *Biotechnol. Bioeng.* **2012**, *109*, 2955–2961. [CrossRef] [PubMed]
- 23. Wójcik, M.; Vázquez Torres, S.; Quax, W.J.; Boersma, Y.L. Sortase mutants with improved protein thermostability and enzymatic activity obtained by consensus design. *Protein Eng. Des. Sel.* **2019**, *32*, 555–564. [CrossRef] [PubMed]
- 24. Wójcik, M.; Szala, K.; van Merkerk, R.; Quax, W.J.; Boersma, Y.L. Engineering the specificity of Streptococcus pyogenes sortase A by loop grafting. *Proteins* **2020**, *88*, 1394–1400. [CrossRef] [PubMed]
- 25. Nikghalb, K.D.; Horvath, N.M.; Prelesnik, J.L.; Banks, O.G.B.; Filipov, P.A.; Row, R.D.; Roark, T.J.; Antos, J.M. Expanding the Scope of Sortase-Mediated Ligations by Using Sortase Homologues. *Chembiochem* **2018**, *19*, 185–195. [CrossRef] [PubMed]
- 26. Schmohl, L.; Bierlmeier, J.; von Kügelgen, N.; Kurz, L.; Reis, P.; Barthels, F.; Mach, P.; Schutkowski, M.; Freund, C.; Schwarzer, D. Identification of sortase substrates by specificity profiling. *Bioorg. Med. Chem.* **2017**, 25, 5002–5007. [CrossRef]
- 27. Malik, A.; Kim, S.B. A comprehensive in silico analysis of sortase superfamily. J. Microbiol. 2019, 57, 431–443. [CrossRef]
- 28. Di Girolamo, S.; Puorger, C.; Castiglione, M.; Vogel, M.; Gébleux, R.; Briendl, M.; Hell, T.; Beerli, R.R.; Grawunder, U.; Lipps, G. Characterization of the housekeeping sortase from the human pathogen Propionibacterium acnes: First investigation of a class F sortase. *Biochem. J.* **2019**, *476*, 665–682. [CrossRef]
- 29. Harms, M.J.; Thornton, J.W. Analyzing protein structure and function using ancestral gene reconstruction. *Curr. Opin. Struct. Biol.* **2010**, 20, 360–366. [CrossRef]
- Pillai, A.S.; Chandler, S.A.; Liu, Y.; Signore, A.V.; Cortez-Romero, C.R.; Benesch, J.L.P.; Laganowsky, A.; Storz, J.F.; Hochberg, G.K.A.; Thornton, J.W. Origin of complexity in haemoglobin evolution. *Nature* 2020, 581, 480–485. [CrossRef]
- 31. Wheeler, L.C.; Lim, S.A.; Marqusee, S.; Harms, M.J. The thermostability and specificity of ancient proteins. *Curr. Opin. Struct. Biol.* **2016**, *38*, 37–43. [CrossRef]
- 32. UniProt Consortium UniProt: A worldwide hub of protein knowledge. *Nucleic Acids Res.* **2019**, 47, D506–D515. [CrossRef] [PubMed]
- 33. Katoh, K.; Rozewicki, J.; Yamada, K.D. MAFFT online service: Multiple sequence alignment, interactive sequence choice and visualization. *Brief Bioinform.* **2017**, 20, 1160–1166. [CrossRef] [PubMed]
- 34. Campen, A.; Williams, R.M.; Brown, C.J.; Meng, J.; Uversky, V.N.; Dunker, A.K. TOP-IDP-scale: A new amino acid scale measuring propensity for intrinsic disorder. *Protein Pept. Lett.* **2008**, *15*, 956–963. [CrossRef] [PubMed]
- 35. Kyte, J.; Doolittle, R.F. A simple method for displaying the hydropathic character of a protein. *J. Mol. Biol.* **1982**, 157, 105–132. [CrossRef]
- 36. Schwartz, G.W.; Zhou, Y.; Petrovic, J.; Fasolino, M.; Xu, L.; Shaffer, S.M.; Pear, W.S.; Vahedi, G.; Faryabi, R.B. TooManyCells identifies and visualizes relationships of single-cell clades. *Nat. Methods* **2020**, *17*, 405–413. [CrossRef] [PubMed]
- 37. Kattke, M.D.; Chan, A.H.; Duong, A.; Sexton, D.L.; Sawaya, M.R.; Cascio, D.; Elliot, M.A.; Clubb, R.T. Crystal Structure of the Streptomyces coelicolor Sortase E1 Transpeptidase Provides Insight into the Binding Mode of the Novel Class E Sorting Signal. *PLoS ONE* **2016**, *11*, e0167763. [CrossRef]
- 38. Bentley, M.L.; Gaweska, H.; Kielec, J.M.; McCafferty, D.G. Engineering the substrate specificity of Staphylococcus aureus Sortase A. The beta6/beta7 loop from SrtB confers NPQTN recognition to SrtA. *J. Biol. Chem.* **2007**, 282, 6571–6581. [CrossRef]
- 39. Zou, Z.; Nöth, M.; Jakob, F.; Schwaneberg, U. Designed Streptococcus pyogenes Sortase A Accepts Branched Amines as Nucleophiles in Sortagging. *Bioconjug. Chem.* **2020**, *31*, 2476–2481. [CrossRef]
- 40. Kruger, R.G.; Dostal, P.; McCafferty, D.G. Development of a high-performance liquid chromatography assay and revision of kinetic parameters for the Staphylococcus aureus sortase transpeptidase SrtA. *Anal. Biochem.* **2004**, *326*, 42–48. [CrossRef]
- 41. Bordoli, L.; Kiefer, F.; Arnold, K.; Benkert, P.; Battey, J.; Schwede, T. Protein structure homology modeling using SWISS-MODEL workspace. *Nat. Protoc.* **2009**, *4*, 1–13. [CrossRef]
- 42. Waterhouse, A.; Bertoni, M.; Bienert, S.; Studer, G.; Tauriello, G.; Gumienny, R.; Heer, F.T.; de Beer, T.A.P.; Rempfer, C.; Bordoli, L.; et al. SWISS-MODEL: Homology modelling of protein structures and complexes. *Nucleic Acids Res.* **2018**, *46*, W296–W303. [CrossRef] [PubMed]
- 43. Arnold, K.; Bordoli, L.; Kopp, J.; Schwede, T. The SWISS-MODEL workspace: A web-based environment for protein structure homology modelling. *Bioinformatics* **2006**, 22, 195–201. [CrossRef] [PubMed]
- 44. Suree, N.; Liew, C.K.; Villareal, V.A.; Thieu, W.; Fadeev, E.A.; Clemens, J.J.; Jung, M.E.; Clubb, R.T. The structure of the Staphylococcus aureus sortase-substrate complex reveals how the universally conserved LPXTG sorting signal is recognized. *J. Biol. Chem.* 2009, 284, 24465–24477. [CrossRef]

45. Varadi, M.; Anyango, S.; Deshpande, M.; Nair, S.; Natassia, C.; Yordanova, G.; Yuan, D.; Stroe, O.; Wood, G.; Laydon, A.; et al. AlphaFold Protein Structure Database: Massively expanding the structural coverage of protein-sequence space with high-accuracy models. *Nucleic Acids Res.* **2022**, *50*, D439–D444. [CrossRef] [PubMed]

- 46. Jumper, J.; Evans, R.; Pritzel, A.; Green, T.; Figurnov, M.; Ronneberger, O.; Tunyasuvunakool, K.; Bates, R.; Žídek, A.; Potapenko, A.; et al. Highly accurate protein structure prediction with AlphaFold. *Nature* **2021**, *596*, 583–589. [CrossRef] [PubMed]
- 47. Bradshaw, W.J.; Davies, A.H.; Chambers, C.J.; Roberts, A.K.; Shone, C.C.; Acharya, K.R. Molecular features of the sortase enzyme family. *FEBS J.* **2015**, *282*, 2097–2114. [CrossRef]
- 48. Comfort, D.; Clubb, R.T. A comparative genome analysis identifies distinct sorting pathways in gram-positive bacteria. *Infect. Immun.* **2004**, *72*, 2710–2722. [CrossRef]
- 49. Dorr, B.M.; Ham, H.O.; An, C.; Chaikof, E.L.; Liu, D.R. Reprogramming the specificity of sortase enzymes. *Proc. Natl. Acad. Sci. USA* **2014**, *111*, 13343–13348. [CrossRef]
- 50. Pedregosa, F.; Varoquaux, G.; Gramfort, A.; Michel, V.; Thirion, B.; Grisel, O.; Blondel, M.; Prettenhofer, P.; Weiss, R.; Dubourg, V.; et al. Scikit-learn: Machine Learning in Python. *J. Mach. Learn Res.* **2011**, *12*, 2825–2830.
- 51. NCBI Resource Coordinators. Database resources of the national center for biotechnology information. *Nucleic Acids Res.* **2017**, 45, D12–D17. [CrossRef]
- 52. Li, W.; Godzik, A. Cd-hit: A fast program for clustering and comparing large sets of protein or nucleotide sequences. *Bioinformatics* **2006**, 22, 1658–1659. [CrossRef] [PubMed]
- 53. Fu, L.; Niu, B.; Zhu, Z.; Wu, S.; Li, W. CD-HIT: Accelerated for clustering the next-generation sequencing data. *Bioinformatics* **2012**, 28, 3150–3152. [CrossRef] [PubMed]
- 54. Altschul, S.F.; Gish, W.; Miller, W.; Myers, E.W.; Lipman, D.J. Basic local alignment search tool. *J. Mol. Biol.* **1990**, 215, 403–410. [CrossRef]
- 55. Edgar, R.C. MUSCLE: Multiple sequence alignment with high accuracy and high throughput. *Nucleic Acids Res.* **2004**, 32, 1792–1797. [CrossRef] [PubMed]
- 56. Guindon, S.; Dufayard, J.-F.; Lefort, V.; Anisimova, M.; Hordijk, W.; Gascuel, O. New algorithms and methods to estimate maximum-likelihood phylogenies: Assessing the performance of PhyML 3.0. *Syst. Biol.* **2010**, *59*, 307–321. [CrossRef] [PubMed]
- 57. Wilkins, M.R.; Gasteiger, E.; Bairoch, A.; Sanchez, J.C.; Williams, K.L.; Appel, R.D.; Hochstrasser, D.F. Protein identification and analysis tools in the ExPASy server. *Methods Mol. Biol.* **1999**, *112*, 531–552.
- 58. Larsson, A. AliView: A fast and lightweight alignment viewer and editor for large datasets. *Bioinformatics* **2014**, *30*, 3276–3278. [CrossRef]
- 59. FigTree. Available online: http://tree.bio.ed.ac.uk/software/figtree/ (accessed on 16 March 2022).