

RESEARCH ARTICLE

PHA synthase variant design using a conditional variational autoencoder

Tuula Tenkanen^{1*}, Anna Ylinen¹, Paula Jouhten², Merja Penttilä¹, Sandra Castillo^{1*}

1 VTT Technical Research Centre of Finland Ltd., Espoo, Finland, **2** Department of Bioproducts and Biosystems, Aalto University, Espoo, Finland

* sandra.castillo@vtt.fi (SC); tuula.tenkanen@vtt.fi (TT)



Abstract

Polyhydroxyalkanoate (PHA) synthases are a group of complex, dimeric enzymes which catalyze polymerization of R-hydroxyacids into PHAs. PHA properties depend on their monomer composition but enzymes found in nature have limited specificities to certain R-hydroxyacids only. In this study, a conditional variational autoencoder was used for the first time to design novel PHA synthases. The model was trained with native protein sequences obtained from Uniprot and was used for the creation of approximately 10 000 new PHA synthase enzymes. Out of these, 16 sequences were selected for *in vivo* validation. The selection criteria included the presence of conserved residues such as catalytic amino acids and amino acids in the dimer interface and structural features like the number of α -helices in the N-terminal part of the enzyme. Two of the 16 novel PHA synthases that had substantial numbers of amino acid substitutions (87 and 98) with respect to the most similar native enzymes were confirmed active and produced poly(hydroxybutyrate) (PHB) when expressed in yeast *S. cerevisiae*. The results show the power of AI based methods to create active variants of highly complex dimer enzymes.

OPEN ACCESS

Citation: Tenkanen T, Ylinen A, Jouhten P, Penttilä M, Castillo S (2026) PHA synthase variant design using a conditional variational autoencoder. *PLoS Comput Biol* 22(3): e1014087. <https://doi.org/10.1371/journal.pcbi.1014087>

Editor: Alexey Onufriev, Virginia Polytechnic Institute and State University, UNITED STATES OF AMERICA

Received: July 4, 2025

Accepted: March 4, 2026

Published: March 19, 2026

Copyright: © 2026 Tenkanen et al. This is an open access article distributed under the terms of the [Creative Commons Attribution License](https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

Data availability statement: Python implementation of the model is available from https://github.com/vttresearch/PHA_cVAE/. Model weights are available on Zenodo at <https://doi.org/10.5281/zenodo.14515368>. Full

Author summary

Enzymes found in nature are limited to the ones that have been beneficial for life during evolution. However, enzymes as proteins whose function arises from their structure are not limited to the ones existing in nature. Therefore, protein design calls for intelligent methods that generate proteins that are expressed, fold, and are active. In this work we developed a deep generative model for PHA synthase variant design. Deep generative models generate new data that resembles the training data. We trained our model using natural polyhydroxyalkanoate (PHA) synthases to generate novel PHA synthase variants. PHA synthases use various monomers to polymerize PHA that has potential as oil-based plastic replacement material. We analyzed the activity of 16 novel PHA synthases we designed and found two of them active. The two active enzymes contained 87 and 98 amino

training and test data is available on Zenodo at <https://zenodo.org/records/17549219>.

Funding: The work was funded by the Jane and Aatos Erkko Foundation (JAES) (<https://jaes.fi/en/frontpage/>) under project 220048, Virtual Laboratory for Biodesign (BIODESIGN) (MP, SC, and TT), and by the Wihuri Foundation (<https://wihurinrahasto.fi/?lang=en>) (Funding for the Center for Young Synbio Scientists and a working grant for T. Tenkanen) (MP and TT). The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

Competing interests: The authors have declared that no competing interests exist.

acid substitutions compared with the closest native PHA synthases. Our work paves the way for the design of novel PHA synthase variants and other enzymes of application interest.

Introduction

Polyhydroxyalkanoates (PHAs) are a group of polymers synthesized by several different bacteria. PHAs have gained attention due to their excellent biodegradability and biocompatibility, making them interesting alternatives to conventional thermoplastic materials. One key research question within PHA field is understanding how the main enzyme, PHA synthase (PhaC), actually carries out the polymerization process and how it can be engineered to be more efficient and accept new monomers [1,2]. Inclusion of new monomers and creation of novel PHA copolymer structures could be an efficient tool for expanding the PHA property space towards new applications.

However, PHA synthase is a challenging enzyme to be modeled and understood with rational methods as the complete crystal structure has been resolved only for one full-length PHA synthase from *Aeromonas caviae* [3]. In addition, PHA synthase acts as a dimer, trimer [4] or multimer [5], containing one or two different subunits. This subunit structure, together with preference for either short-chain-length (scl) or medium-chain-length (mcl) monomers, defines the class of a certain PHA synthase. Class I, III and IV PHA synthases prefer scl-monomers with only three to five carbons, while class II PHA synthases prefer mcl-monomers containing six to 14 carbons. Class I and II PHA synthases act as homodimers containing two similar PhaC subunits although recently Assefa *et al.* [4] suggested that class I PHA synthase from *Brevundimonas* sp. KH11J01 is active as a trimer. In contrast, class III and IV PHA synthases have in addition to PhaC subunit a PhaE or a PhaR subunit, respectively [6] and exist as heterodimers or heteromultimers [5].

Instead of rational design, computational tools can be used. Computational protein design tools such as FuncLib [7] and CADENZ [8] can be used to generate libraries with enzyme variants. Enzymes with enhanced activity or even new substrate specificities have been found by screening these libraries [9]. In addition, to these tools generative AI models, such as variational autoencoders (VAEs) or large language models (LLMs), can be used to generate libraries with novel enzymes. While VAEs have been successfully applied to design novel proteins such as metalloproteins, luciferase enzymes, and simpler proteins such as human-like phenylalanine hydroxylases [10–12], their potential for creating more complex enzymes such as dimers, like PHA synthases, has been understudied.

In this work, we designed two novel (i.e., AI generated) class I PHA synthase sequences and demonstrated their activity by polymerization of 3-hydroxybutyrate (3HB) monomers in our previously developed yeast *Saccharomyces cerevisiae* strain expressing a 3HB-CoA synthesizing pathway [13]. New PHA synthase sequences were designed using a generative AI model, specifically a conditional variational auto-encoder (cVAE) PHA_cVAE.

Results

Choice of a deep learning model

We trained an autoencoder model that takes as input the structural and sequence-based features of enzymes and outputs the amino acid sequences of enzymes (Fig 1). The model also receives a condition vector that represents the enzyme class. We used a dataset of PHA synthases, lipases, and partial PHA synthase sequences for the model training. We found that using bidirectional LSTM layers and multi-head self-attention blocks improved the model performance.

We defined the loss function of the model as a combination of the reconstruction loss and the Kullback-Leibler divergence (KL divergence) [14], which measures how much the latent space differs from the standard normal distribution. The KL divergence term can be scaled by a factor called beta that controls the trade-off between reconstruction and disentanglement of the latent features [15]. However, a high beta value can lead to a posterior collapse [16], a phenomenon where the latent variables become independent of the input and the decoder ignores them. To mitigate this problem, we trained the model in stages, gradually increasing the beta parameter from 0.01 to 1.

Model performance was evaluated by analyzing the quality of the produced PHA synthase sequences. Class I PHA synthase sequences generated with PHA_cVAE were compared with class I PHA synthase sequences produced with a

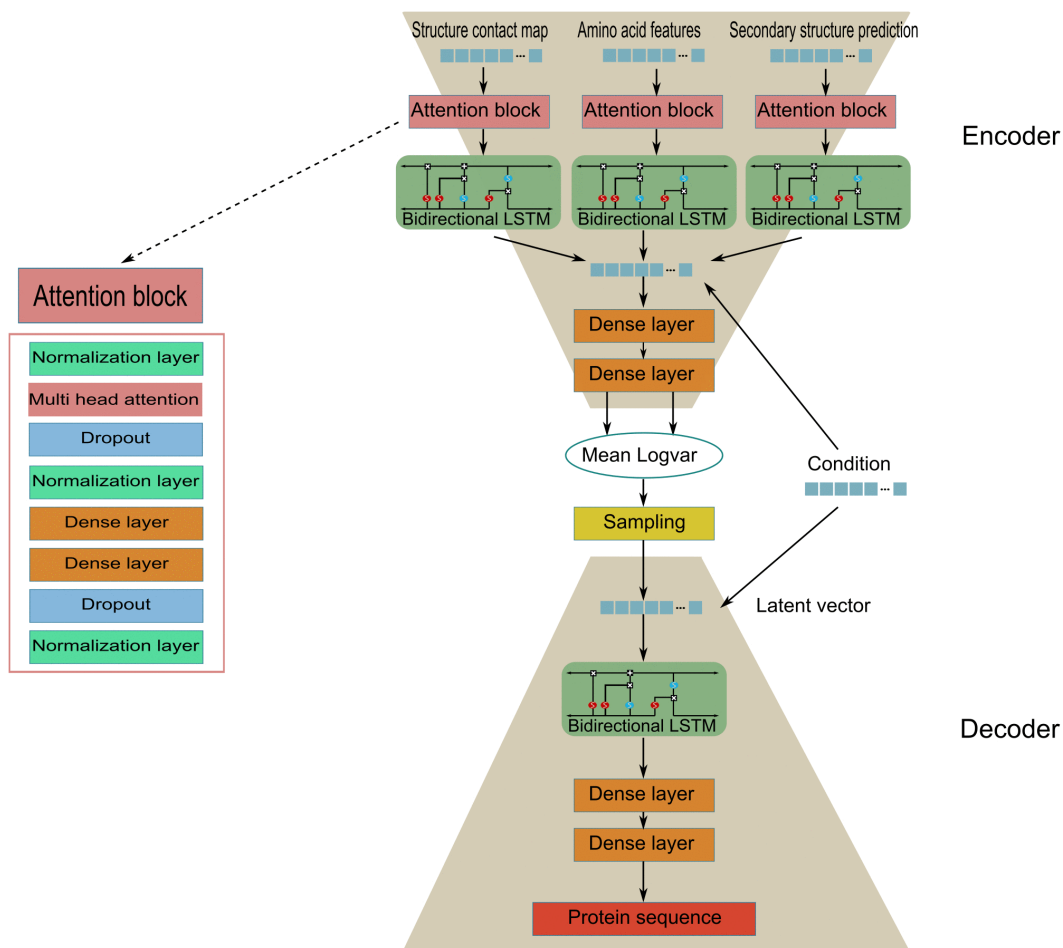


Fig 1. Scheme of the deep learning model used to generate novel PHA synthases.

<https://doi.org/10.1371/journal.pcbi.1014087.g001>

fine-tuned large language model (LLM) ProGen2-small [17] as well as native class I PHA synthase sequences with 100 random mutations. The evaluation focused on class I PHA synthase sequences as we had decided to generate class I PHA synthase sequences with the model. In addition, we generated sequences with PHA_cVAE in all four PHA synthase classes and evaluated them using three different criteria: coherence of BLAST results with the generated class, confidence of the model to terminate the sequence and amino acid composition.

Some sequences generated with PHA_cVAE contained several long homorepeats and did not thus seem to be proper PHA synthase sequences (only 0.3% of the input sequences contained at least one homorepeat longer than 10 amino acids). Therefore, we first evaluated if the generated sequences had also rare amount of some amino acids. We analyzed the percentage range of each amino acid in the native sequences (e.g., amount of alanine residues was 4–16% in the native sequences) and compared if the generated sequences had similar percentage of each amino acid as the native sequences (Table 1). We then used this amino acid composition as a filter before further evaluation to have only meaningful sequences in the analysis.

Next, we predicted structures of the generated class I PHA synthase variants with Boltz-1 [18] and calculated the average pLDDT score [19] from the predicted structures. Furthermore, we aligned the predicted structures with PHA synthase from *Chromobacterium* sp. USM2 and calculated average TM-score [20]. The average pLDDT score was slightly higher for the PHA_cVAE generated sequences than for the ProGen2-small generated sequences (Table 1). Therefore, the structure prediction model Boltz-1 was slightly more confidently predicting the structures of the PHA synthase variants generated by the PHA_cVAE model. Furthermore, pLDDT scores of sequences generated with both generative models were higher than pLDDT scores of native PHA synthase sequences with 100 random mutations (Table 1). The average TM-scores were similar for PHA synthase variants generated by both models as well as for native PHA synthase sequences with random mutations (Table 1). Distribution plots for pLDDT scores and TM-scores are presented in S1 Fig.

We also calculated the proportion of generated sequences that contained the three catalytic amino acids (Cys, Asp, His) in the right positions. Both PHA_cVAE and ProGen2-small mainly generated sequences that contained catalytic amino acids at right positions, but ProGen2-small was slightly better at it (Table 1). The ProGen2-small was also generating more variability in the sequences. Clustering the sequences with CD-HIT [21] to 0.95 of similarity produced more clusters for ProGen2-small generated sequences than for PHA_cVAE generated sequences (Table 1).

We then analyzed if PHA_cVAE had learned to distinguish the different PHA synthase classes from each other. We generated sequences in all four PHA synthase classes and BLASTed [22] the generated sequences against native PHA synthase sequences. We then evaluated if the hit with the lowest e-value was from the correct class. To compare the two

Table 1. Evaluation of model performance^a.

Class of generated PHA synthase variants	PHA_cVAE				ProGen2-small	Random mutations
	I	II	III	IV	I	I
Average pLDDT score	88.3	–	–	–	85.4	78.4
Average TM-score	0.90	–	–	–	0.89	0.88
Correct catalytic amino acid triad	96%	–	–	–	97.6%	–
Amount of clusters with 0.95 similarity threshold	160	–	–	–	544	–
Amino acid composition	33%	37%	39%	28%	77%	–
Coherence of BLAST results with the class generated	100%	100%	93%	97%	96%	–
Confidence of the model to terminate the sequence	72%	63%	34%	66%	–	–

^aPHA synthase variants were generated in all four PHA synthase classes with PHA_cVAE and in class I with ProGen2-small. Furthermore, native class I PHA synthase sequences were randomly mutated, introducing 100 mutations into each sequence. “–” denotes no analysis.

<https://doi.org/10.1371/journal.pcbi.1014087.t001>

models we also applied the same analysis for ProGen2-small generated sequences. All class I PHA synthase sequences generated with PHA_cVAE were from correct class while 4% of sequences generated with ProGen2-small were not evaluated as class I PHA synthase variants (Table 1). Thus, PHA_cVAE was slightly better in producing sequences of correct PHA synthase class.

Furthermore, we evaluated how confident PHA_cVAE was in terminating the sequences. The analysis was conducted prior to filtering the sequences based on amino acid composition. Since all inputs were required to be 700 amino acids in length, shorter sequences were padded to achieve uniform length. Consequently, the generated sequences also included padding at the end to indicate termination. In some cases, the model introduced termination indicators within the middle of sequences, suggesting uncertainty in determining the correct endpoint (Table 1).

Selection of 16 novel PHA synthase variants

We used PHA_cVAE to create novel class I PHA synthase sequences. We focused on class I as it is the most studied PHA synthase class. Furthermore, we wanted to focus on a PHA synthase class that is active as a homodimer (i.e., either class I or II). Generation of class III or IV enzymes would bring additional dimension as in addition to PhaC also PhaE or PhaR subunits are needed, respectively. More than 10 000 class I PHA synthase sequences were generated to have variance in the sequences. Only sequences with similar amino acid composition as native class I PHA synthases were saved. As it is currently challenging, if not impossible, to experimentally characterize 10 000 different enzyme sequences, we carried out several filtration steps to select 16 most interesting sequences for an *in vivo* activity test. First we removed duplicates and sequences without correct active site triad (Cys, Asp, His). This first filtering step removed 10% of the initial sequences. Next we clustered the sequences based on similarity with CD-HIT [21] and selected one sequence from each cluster leaving 715 sequences. We then predicted the structures for these 715 sequences with AlphaFold [19], compared them with PHA synthase from *Chromobacterium* sp. USM2 and analyzed if the generated sequences contained 20 different amino acids that were shown to be conserved among native PHA synthases from four different PHA synthase classes [23] and an arginine residue (Arg365 in *Chromobacterium* sp. USM2) that is conserved in classes I and II. In addition we evaluated if the sequences contained hydrophobic amino acids at the dimer interface obtained from Chek *et al.* [23]. Next, we evaluated the amount of α -helices in the N-terminal part as Kim *et al.* [24] analyzed that five α -helices in the N-terminal are required for PhaC_{1_{cn}} to function properly. In addition, we analyzed the length of the sequences, presence of termination indicators inside the sequences and similarity with the closest native enzymes. We then selected 42 sequences from the 715 sequences based on the information gathered above. We selected mostly sequences that contained all conserved amino acids, but selected also some sequences with some variability in these. Length of the selected sequences were 570–621 aa, sequence identities with closest native enzyme were 50–92%, and they all contained 5 α -helices in the N-terminal part. Next, we evaluated tunnels to the active site cavity with CAVER [25] and selected 16 sequences for wet lab validation. All of the selected 16 sequences had two tunnels to the active site cavity with bottleneck radius of at least 1.3 Å. Furthermore, all of the selected sequences apart from PhaC_{VAE6} and PhaC_{VAE10} contained all conserved amino acids and hydrophobic amino acids in the dimer interface.

Fifteen of the resulting 16 selected enzymes (Table 2) showed highest sequence identity to PHA synthases either from *Legionella* or *Janthinobacterium* species. To assess whether this bias was a result of the selection process or if all generated sequences were similar to these species, we constructed a phylogenetic tree from the initially generated sequences by clustering all the generated sequences with CD-hit using 50% similarity threshold and including one representative sequence from each cluster. In addition, all of the sequences selected for *in vivo* activity tests were added to the phylogenetic tree (Fig 2). The sequences generated using the model had remarkably higher variability than the sequences selected for *in vivo* tests showing the model ability to avoid the usual issue of VAEs that reduce variance in the sampling [26]. Instead, using conserved residues as a selection criterion caused some bias in the selection process.

Table 2. PHA_cVAE generated PHA synthase enzymes^a.

Enzyme ID	Active	Uniprot ID of closest native specie	Closest native specie	Sequence identity (%)	Changes	Rare changes	Novel change combinations
PhaC _{VAE1}	Active	A0A0W0ZA05	<i>Legionella shakespearei</i>	85	87	7 (0)	52
PhaC _{VAE2}	Active	A0A0W0RV18	<i>Legionella bozemanee</i>	83	98	2 (0)	51
PhaC _{VAE3}	Inactive	A0A2N6IF07	<i>Janthinobacterium sp. ROICE36</i>	88	82	17 (2)	1260
PhaC _{VAE4}	Inactive	A0A5C4NWS4	<i>Janthinobacterium lividum</i>	88	70	14 (0)	801
PhaC _{VAE5}	Inactive	A0A1I9XWC6	<i>Janthinobacterium sp. 1_2014MBL_MicDiv</i>	85	88	15 (1)	835
PhaC _{VAE6}	Inactive	A0A238KKZ1	<i>Actibacterium lipolyticum</i>	61	237	65 (9)	15907
PhaC _{VAE7}	Inactive	A0A2N0HPF3	<i>Janthinobacterium sp. 64</i>	84	92	20 (8)	1918
PhaC _{VAE8}	Inactive	A0A5C4NWS4	<i>Janthinobacterium lividum</i>	89	68	16 (3)	818
PhaC _{VAE9}	Inactive	A0A377RVJ7	<i>Janthinobacterium lividum</i>	90	71	7 (0)	434
PhaC _{VAE10}	Inactive	A0A0W0ZA05	<i>Legionella shakespearei</i>	73	157	30 (9)	4695
PhaC _{VAE11}	Inactive	A0A5C4NWS4	<i>Janthinobacterium lividum</i>	90	56	14 (0)	759
PhaC _{VAE12}	Inactive	A0A0W0ZA05	<i>Legionella shakespearei</i>	76	137	16 (0)	1113
PhaC _{VAE13}	Inactive	A0A5C4NWS4	<i>Janthinobacterium lividum</i>	91	53	8 (0)	438
PhaC _{VAE14}	Inactive	A0A5C4NWS4	<i>Janthinobacterium lividum</i>	88	76	14 (2)	745
PhaC _{VAE15}	Inactive	A0A5C4NWS4	<i>Janthinobacterium lividum</i>	87	91	15(1)	1215
PhaC _{VAE16}	Inactive	A0A378IN02	<i>Legionella cinцинnatiensis</i>	85	89	5 (0)	92

^aThe “Active” column shows if the PHA synthase variant showed activity for 3HB-CoA polymerization when expressed *in vivo* in yeast *S. cerevisiae*. The second column is the Uniprot ID of the most similar native PHA synthase and the third column is the organism where the most similar native PHA synthase was found. The “Sequence identity” column shows percentage of identical positions with the most similar native PHA synthase. The “Changes” column shows the number of amino acid differences between the generated PHA synthase and the most similar native PHA synthase. The “Rare changes” column shows the number or amino acid positions where the amino acid present in the novel PHA synthase is present in less than 100 native class I PHA synthases in the same position. Inside the parenthesis we show the number of novel changes (i.e., amino acid positions where amino acid present in the novel enzyme is not found in any native class I PHA synthase). The column “Novel change combinations” shows the number of two change combinations that were not present in any native class I PHA synthase.

<https://doi.org/10.1371/journal.pcbi.1014087.t002>

To further evaluate how the selected PHA synthases were distributed in the latent space we visualized the selected PHA synthases together with 1000 randomly selected PHA_cVAE generated sequences and 1674 randomly selected sequences from the train and test dataset (S1 File). As expected, those generated sequences that were similar to PHA synthases from *Legionella* sp., based on sequence similarity, were close to native PHA synthases from *Legionella* sp. in the latent space. Generated sequences similar to *Janthinobacterium* sp. formed two clusters with one cluster containing sequences PhaC_{VAE3} and PhaC_{VAE4} and other cluster containing the rest of the generated sequences that were similar to *Janthinobacterium* sp., including PhaC_{VAE7} and native PHA synthases from *Janthinobacterium* sp. Although the closest native PHA synthase for PhaC_{VAE3}, PhaC_{VAE4} and PhaC_{VAE7} was from *Janthinobacterium* sp. the sequence similarity between both PhaC_{VAE3} and PhaC_{VAE7} and PhaC_{VAE4} and PhaC_{VAE7} is approximately 70%. The difference in the sequence similarity might therefore be one reason why these sequences are not next to each other in the latent space. However, another reason might be that the latent space distributes the data not only based on amino acid sequence, but captures also other features of the enzymes. Therefore, the structured latent space can be used to find relationships between enzymes that can not be seen with sequence alignment.

In vivo activity test of the 16 novel PHA synthases

Together with PhaC_{LS}, PhaC_{J1}, and PhaC_{1Cn} encoding genes, the genes encoding for the selected 16 PHA_cVAE generated novel PHA synthases were individually integrated into chromosome X of *S. cerevisiae*, more specifically into X-4 EasyClone loci [27]. Parent strain contained three copies of the 3-hydroxybutyryl-CoA (3HB-CoA) pathway, including

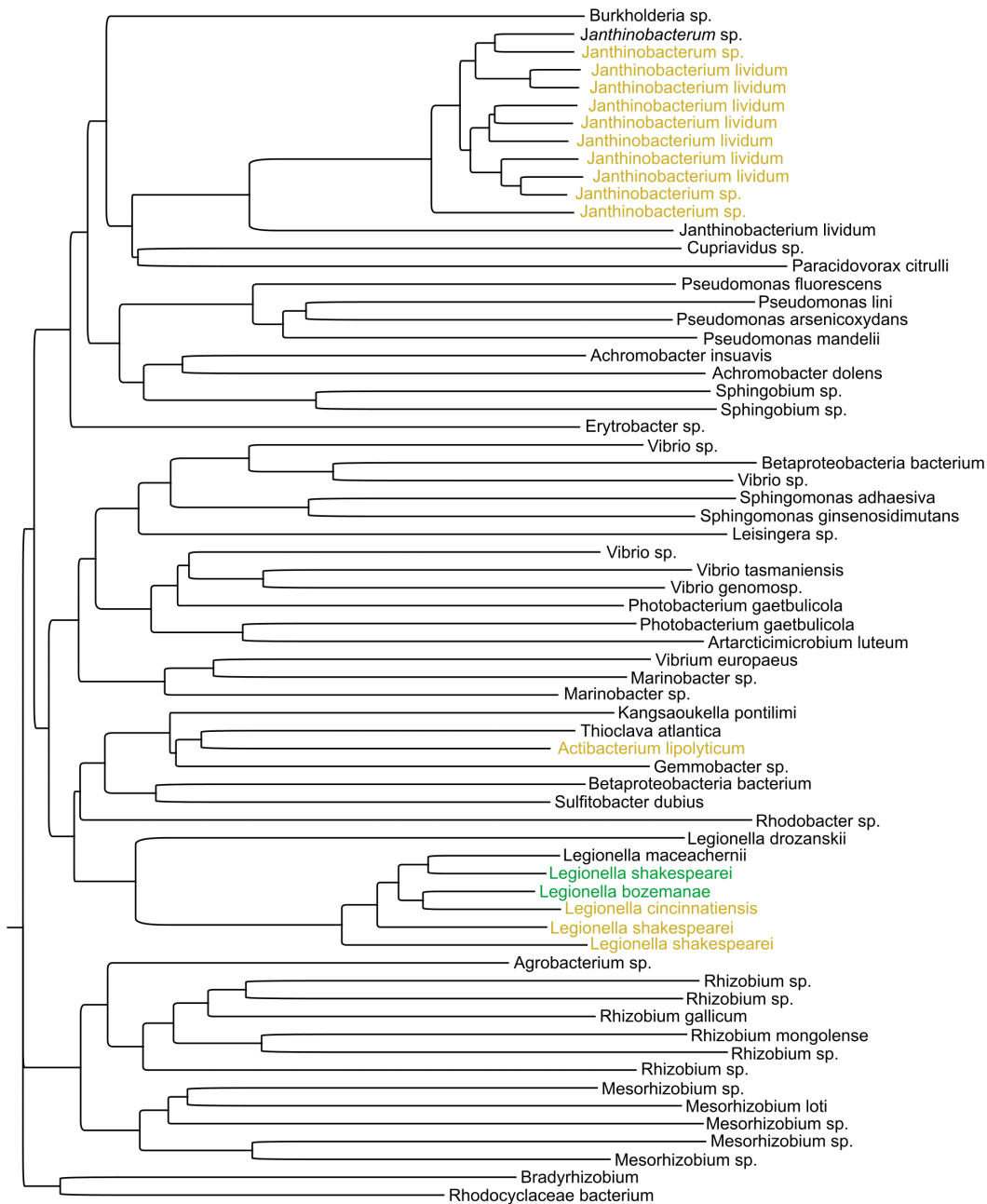


Fig 2. Phylogenetic tree of PHA_cVAE generated PHA synthase sequences. The generated sequences were clustered and one representative from each cluster is included in the tree. In addition, sequences selected for *in vivo* activity analysis were added. Active sequences are colored with green and non active with orange. Labels in the tree describe the species of the closest native PHA synthase.

<https://doi.org/10.1371/journal.pcbi.1014087.g002>

acetyl-CoA acetyltransferase (PhaA) and acetoacetyl-CoA reductase (PhaB1) (Fig 3A) (Table 3). The activities of the enzymes were first assessed by staining with Nile red that binds on the surface of the PHA granules [28]. The Nile red staining suggested that two of the novel PHA synthases were active and polymerized 3HB monomers (Fig 3B). The strains expressing PHAC_{VAE1} and PHAC_{VAE2} showed significantly higher based on two-tailed paired Student's t-Test

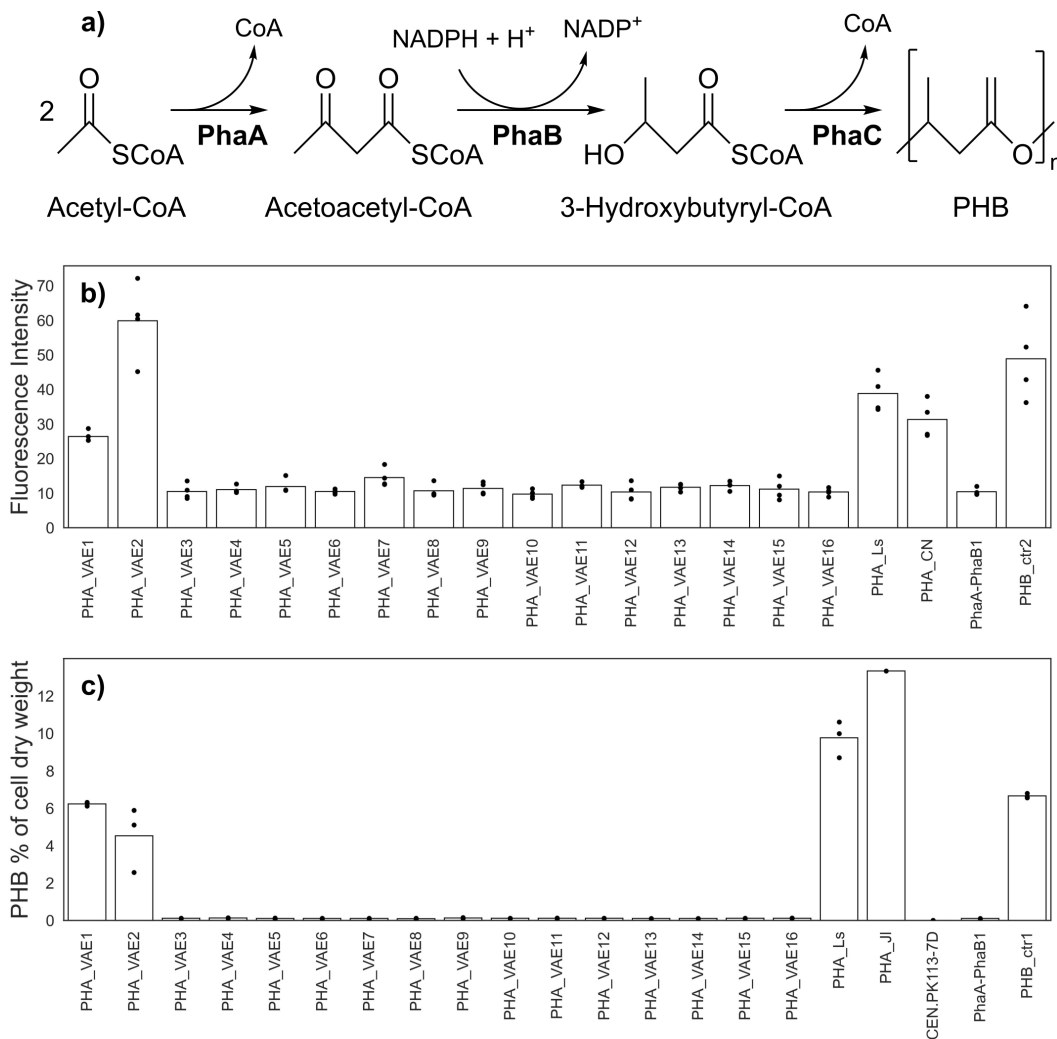


Fig 3. Biosynthesis pathway of PHB and results of PHB accumulation in vivo. (A) PHB biosynthesis pathway by PhaA, PhaB1 and PhaC. **(B)** Fluorescence intensities of strains grown on 24-well plate and stained with Nile red method. Four technical replicates were used for each strain. Bars show the average value of technical replicates, and circles show the values of individual measurements. **(C)** The amount of PHB as % of cell dry weight quantified with GC-MS from strains grown in shake flasks. Three biological replicates were used for strains PHA_VAE1-PHA_VAE16 and PHA_Ls. For control strains CEN.PK113-7D, PhaA-PhaB1, and PHB_ctr1 same biological replicate was cultured in three different shake flasks. For strain PHA_JI, one biological replicate was cultivated. Bars show the average value of replicate cultures, and circles show the value of each replicate culture.

<https://doi.org/10.1371/journal.pcbi.1014087.g003>

($p < 0.003$) fluorescence (26.4 ± 1.6 Relative fluorescence units (RFU) and 59.8 ± 11.1 RFU, respectively) than the PHA negative control strain PhaA-PhaB1 (10.4 ± 1.0 RFU). The nonzero fluorescence intensity of the negative control is explained by the binding of Nile red to other intracellular structures such as lipids droplets [29]. For comparison, positive controls PHA_Ls, PHA_Cn, and PHB_ctr2 showed similar fluorescence intensities (38.8 ± 5.4 RFU, 31.3 ± 5.4 RFU, and 48.8 ± 12.1 RFU, respectively) as the strains expressing the novel PHA_{VAE1} and PHA_{VAE2}.

Next, the activities of the 16 novel PHA synthases were assessed by growing the yeast strains in shake flasks for 72 h and analyzing PHB content from the lyophilized biomass with precise gas chromatography mass spectrometry (GC-MS) method. The analysis of the PHB content in the biomass confirmed the earlier observations from the Nile red staining (Fig 3C). Strains expressing the novel PHA synthases, PHA_VAE1 and PHA_VAE2, accumulated $6.2 \pm 0.1\%$ and

Table 3. Strains and enzymes used in this study.

Strains			
Name	Description	Code	References
CEN.PK113-7D	<i>S. cerevisiae</i> (MATa HIS3 URA3 LEU2 TRP1 MAL2-8c SUC2)	H3887	a
PHB_ctr1	H3887 with integration of <i>PhaA</i> , <i>PhaB1</i> and <i>PhaC1_{Cn}</i> genes into X-3 EasyClone locus	H5696	[30]
PHB_ctr2	H3887 with integration of <i>PhaA</i> , <i>PhaB1</i> and <i>PhaC1_{Cn}</i> genes into X-4, XII-5 and XI-3 EasyClone loci	H5697	This article
PhaA-PhaB1	H3887 with integration of <i>PhaA</i> and <i>PhaB1</i> genes into X-3, XI-2 and XII-2 EasyClone loci	H6135	This article
PHA_VAE1	H6135 with integration of <i>phaC_{VAE1}</i> into X-4 EasyClone locus	H6772	This article
PHA_VAE2	H6135 with integration of <i>phaC_{VAE2}</i> into X-4 EasyClone locus	H6773	This article
PHA_VAE3	H6135 with integration of <i>phaC_{VAE3}</i> into X-4 EasyClone locus		This article
PHA_VAE4	H6135 with integration of <i>phaC_{VAE4}</i> into X-4 EasyClone locus		This article
PHA_VAE5	H6135 with integration of <i>phaC_{VAE5}</i> into X-4 EasyClone locus		This article
PHA_VAE6	H6135 with integration of <i>phaC_{VAE6}</i> into X-4 EasyClone locus		This article
PHA_VAE7	H6135 with integration of <i>phaC_{VAE7}</i> into X-4 EasyClone locus		This article
PHA_VAE8	H6135 with integration of <i>phaC_{VAE8}</i> into X-4 EasyClone locus		This article
PHA_VAE9	H6135 with integration of <i>phaC_{VAE9}</i> into X-4 EasyClone locus		This article
PHA_VAE10	H6135 with integration of <i>phaC_{VAE10}</i> into X-4 EasyClone locus		This article
PHA_VAE11	H6135 with integration of <i>phaC_{VAE11}</i> into X-4 EasyClone locus		This article
PHA_VAE12	H6135 with integration of <i>phaC_{VAE12}</i> into X-4 EasyClone locus		This article
PHA_VAE13	H6135 with integration of <i>phaC_{VAE13}</i> into X-4 EasyClone locus		This article
PHA_VAE14	H6135 with integration of <i>phaC_{VAE14}</i> into X-4 EasyClone locus		This article
PHA_VAE15	H6135 with integration of <i>phaC_{VAE15}</i> into X-4 EasyClone locus		This article
PHA_VAE16	H6135 with integration of <i>phaC_{VAE16}</i> into X-4 EasyClone locus		This article
PHA_Ls	H6135 with integration of <i>phaC_{Ls}</i> into X-4 EasyClone locus	H6774	This article
PHA_Jl	H6135 with integration of <i>phaC_{Jl}</i> into X-4 EasyClone locus	H6775	This article
PHA_Cn	H6135 with integration of <i>phaC1_{Cn}</i> into X-4 EasyClone locus	H6776	This article
Enzymes			
Name	Description		References
PhaA	Acetyl-CoA acetyltransferase from <i>Cupriavidus necator</i> , GenBank KP681582		[31]
PhaB1	Acetoacetyl-CoA reductase from <i>C. necator</i> , GenBank KP681583		[31]
PhaC1 _{Cn}	PHA synthase from <i>Cupriavidus necator</i> , GenBank KP681584		[31]
PhaC _{Ls}	PHA synthase from <i>Legionella shakespearei</i> DSM 23087, Uniprot A0A0W0ZA05		This article
PhaC _{Jl}	PHA synthase from <i>Janthinobacterium lividum</i> , Uniprot A0A5C4NWS4		This article
PhaC _{VAE1} -phaC _{VAE16}	Novel PHA synthases generated by VAE model		This article

*Strain was kindly provided by Dr. P. Kötter (Institut für Mikrobiologie, J.W. Goethe Universität Frankfurt, Germany).

<https://doi.org/10.1371/journal.pcbi.1014087.t003>

4.5 ± 1.7% of PHA as % of CDW, respectively, demonstrating activity of the novel PHA synthases. For comparison, PHB content in the biomass of the CEN.PK113-7D was below the detection limit and the strain expressing only PhaA-PhaB1 and the strains expressing the other novel PHA synthases (PHA_VAE3-PHA_VAE16) showed only a trace amount of 0.05-0.15% of PHB in CDW. The trace amount results likely from methanolysis of an excess of non-polymerized 3-HB-CoA formed by PhaB1. Strains with native PHA synthases from, *C. necator* (PHB_ctr1), *L. shakespearei* (PHA_Ls) and *J. lividum* (PHA_Jl) accumulated 6.7 ± 0.1%, 9.8 ± 1.0%, and 13.3% PHB of CDW, respectively. As PHA titers are linked to amount of formed biomass, the cell growth was followed at 0h, 24h, and 72h [S2 Fig](#). During the first 24 h, almost all strains reached their highest OD₆₀₀ the only exception being strain PHB_ctr1 which showed minor OD₆₀₀ increase after 24 h.

Sequence analysis of the 16 novel PHA synthases

To identify the features that distinguished the active and inactive novel PHA synthases, we performed a comparative analysis of the sequences by structurally aligning them with their closest native PHA synthases. Most of the 16 enzymes tested experimentally were similar to the native PHA synthases from two species: *Legionella* sp. and *Janthinobacterium* sp. (see [Table 2](#)). Two enzymes showed activity *in vivo* with the selected substrate, and they were similar to the native PHA synthases from *Legionella shakespearei* and *Legionella bozemanæ*. PhaC_{VAE1} and PhaC_{VAE2} had 87 and 98 amino acid substitutions, respectively, compared to their closest native PHA synthase. The exact location of the substitutions in PhaC_{VAE2} is visualized in [Fig 4](#). The other enzymes, which were not active *in vivo*, had 53–237 substitutions compared to the most similar native PHA synthase. Despite substantial number of amino acid substitutions the structural motifs of the generated enzymes have remained near identical when compared with the closest native PHA synthase ([Fig 4](#)). In PhaC_{VAE2} the only differences found were shorter β 3- and β 4-strands ([S3 Fig A](#)) and two missing short helices ([S3 Fig B](#)). Also, PhaC_{VAE6} with 237 substitutions have preserved majority of the structural motifs found in the closest native PHA synthase ([S3 Fig C](#)).

Next, we evaluated the frequency and novelty of the amino acid substitutions in each enzyme. The active enzymes exhibited less rare substitutions (amino acids with a frequency of less than 100 in native PHA synthases at the same position) and no novel substitutions (amino acids absent in native PHA synthases at the same position). The mean number of rare substitutions for the active enzymes was 4.5, while for the inactive enzymes it was 18.3. We also analyzed the pairwise combinations of substitutions in each enzyme. The active enzymes had 52 and 51 novel combinations, respectively, while the inactive enzymes had from 92 to 15907 novel combinations (see [Table 2](#)).

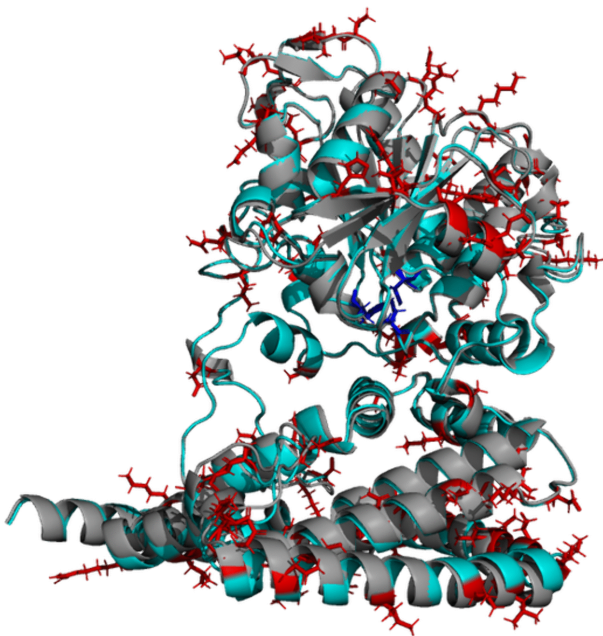


Fig 4. Comparison of amino acid substitutions in PhaC_{VAE2}. AlphaFold structure of the active PhaC_{VAE2} (in cyan) aligned with AlphaFold structure of the closest native PHA synthase (in grey). The 98 different amino acids with respect to the most similar native enzyme are marked with red. Catalytic amino acids are shown with blue.

<https://doi.org/10.1371/journal.pcbi.1014087.g004>

Discussion

In this study we created 16 novel PHA synthases with a conditional variational autoencoder. We found that two of these could effectively polymerize 3HB-CoA into PHB with accumulation levels similar to previously obtained by Ylinen *et al.* when expressing *C. necator* PHA synthase [30]. Similar levels have previously been reported also for other engineered *S. cerevisiae* strains carrying the PHB pathway from *C. necator* (5.2-9% of CDW) [29,32,33], indicating that our novel PHA synthases were as efficient in polymerizing 3HB-CoA as the PHA synthase from *C. necator* when expressed in yeast. However, since it is difficult to normalize the effect of other parameters on PHA accumulation *in vivo*, precise comparison of different enzymes is difficult. For example, monomer availability has great impact on PHA accumulation, as shown in studies where replacement of PhaB1 with an NADH dependent alternative and the use of xylose or cellobiose as the carbon source instead of glucose boosted the carbon flux towards 3HB-CoA and result in higher PHB accumulation levels of 14–21% of CDW in *S. cerevisiae* [30,34,35]. Despite this, the *in vivo* approach is a rather easy method for assessing the capability of a PHA synthase to polymerize a certain monomer in relevant conditions. In future, activity screening could be made by using Nile red staining [28] as Nile red staining results corresponded well with GC-MS analysis. Alternatively activity screening could be made by cultivating microbes containing the synthases in a 96 well plate and analyzing PHA accumulation using Fourier transform infrared spectroscopy (FTIR) [36]. In either option, precursors of desired PHA monomers could be fed to the cultures to evaluate the ability of the synthases to use desired monomers as their substrates [28,37,38]. Finally, GC-MS can be used to confirm the monomer composition of PHA polymerized by active PHA synthase variants.

In addition, the selection process before wet lab activity screening can be improved in future by taking into account the number and type of amino acid substitutions with respect to the most similar native PHA synthase sequence. The two active novel PHA synthase variants contained significantly lower number of novel combinations of substitutions than the novel inactive PHA synthase variants, suggesting that some combinations may be detrimental for the enzyme activity. Furthermore, using the conserved residues as one selection criterion should be reconsidered as it caused some bias during the selection process. Although the residues used were shown to be conserved in some native PHA synthases from all four PHA synthase classes [23] they were not conserved in all of the sequences used for training the model. Therefore, the model also suggested alternative residues for some of these positions. Although majority of the sequences used for training are not validated as active PHA synthases, it is possible that using this selection criterion removed potential active enzymes from wet lab validation. In addition, at the time of the study C-terminal catalytic domains of two PHA synthases from *Cupriavidus necator* [39,40] and *Chromobacterium sp. USM2* [23] were available, but crystal structure of a full-length PHA synthase was lacking. The recently published full-length crystal structure of PhaC from *A. caviae* (PhaC_{Ac}) gives new insight into the dimerization of the enzyme. According to the new crystal structure there are direct contacts between N-terminal residues from both protomers. Therefore, the new knowledge do not support using the hydrophobic residues that were earlier suggested to form contact area between two protomers [23] as one filter. However, with the new crystal structure we could also analyze how the tunnels obtained with CAVER align with the putative egress tunnel presented by Chek *et al.* [3]. We aligned PhaC_{VAE1} and the tunnels obtained with CAVER with triethylene glycol (TEG)-bound form of the crystal structure of PhaC_{Ac} and found that TEG-molecules were aligning with one of the two tunnels (S4 Fig A). In addition, when we docked 3HB-CoA to PhaC_{VAE1} with Boltz-1 the docked 3HB-CoA went through the second tunnel obtained with CAVER (S4 Fig B). Therefore, the tunnels obtained with CAVER seem to be putative substrate entry and product egress tunnels which supports using CAVER in the selection process also in the future.

In future, training dataset and protein representation of the input sequences could also be further improved. Generative AI models typically require large amount of data for effective training. In addition to PHA synthase sequences we added lipase sequences to the training dataset as lipases share some structural similarities with PHA synthases [23]. The effect of adding lipases to the training dataset on the quality of the generative model is unclear. It is possible that some of the lipases have divergent structures that could interfere with the model's ability to learn the features of PHA synthases.

However, the number of PHA synthase sequences in public databases is increasing rapidly. Therefore, in the future the model could be trained exclusively with PHA synthase sequences. In addition, we utilized in this work contact maps predicted by trRosetta [41] to represent the 3D structures of the input sequences for the deep learning model. Contact maps predicted by trRosetta contain information about the distances and orientation of amino acid residues that are in contact with each other in the enzyme structure. Thus, these contact maps describe enzymes better than their protein sequences. Predicting contact maps with trRosetta is fast and can be used easily on a data set of thousands of enzymes. However, the accuracy and quality of the protein representation, and thus also model's performance, could be enhanced in future by the new advances in protein structure prediction such as AlphaFold [19] or RoseTTAFold [42].

Finally, VAEs are efficient generative models which have been used successfully in the past to create novel proteins [10–12]. However, training of VAEs can be difficult and unstable due to various challenges during the training process, such as posterior collapse, vanishing gradients, over-fitting, and dataset bias [43]. We addressed the problem of posterior collapse by adopting a phase training approach increasing gradually the beta parameter in the loss function. Problems like posterior collapse, loss of variability or bias could also be avoided by using alternative types of generative models, such as diffusion models [44] or large language model (LLM) [17], but these models have other disadvantages like large training costs of LLM and no possibility for smart sampling from latent distribution. To compare the quality of sequences generated by PHA_cVAE with those generated by a fine-tuned LLM (ProGen2-small [17]) we predicted structures of the generated sequences and calculated average pLDDT score. In addition, we aligned the structures with a native PHA synthase and calculated average TM-score. The average pLDDT score and TM-score are slightly better for the PHA_cVAE generated sequences. However, the ProGen2-small is able to generate more diversity than PHA_cVAE, based on sequence clustering. In addition, amino acid composition was better in ProGen2-small generated sequences, but PHA_cVAE was slightly better in generating sequences of correct PHA synthase class after filtering the sequences according to amino acid composition. Independently of the quality of the the generated PHA synthases, the main reason for using cVAE instead of LLM is the possibility of smarter sampling from the latent distribution, as it can be highly customized. Sevgen *et al.* [10] designed novel phenylalanine hydroxylases (PAH) by sampling sequences around human PAH from the latent distribution and some of the novel sequences showed increased activity when compared with the human PAH. In the future, we could sample around the sequences that showed activity and try to find improved versions of them. Smart sampling could also be used to design PHA synthases that can utilize a wider or more specific range of substrates. For instance, sampling sequences from latent distribution between a class I native PHA synthase and a class II native PHA synthase could allow generation of novel PHA synthases with substrate specificity towards both scl- and mcl-monomers. In addition, smart sampling is not the only benefit of the structured latent space. As the regularized latent space groups similar inputs closely it enables relationship discovery of enzymes beyond sequence comparison.

Conclusion

In this study, a conditional variational autoencoder was used for the first time to create novel PHA synthases. Despite 87 and 98 amino acid substitutions in comparison to the closest native PHA synthases the two PHA synthases were active and produced PHB in yeast *S. cerevisiae*. Ultimately these, or other novel sequences designed in future, could expand possibilities to polymerize different PHA monomers and adjust PHA material properties into new application areas.

The success rate of our low-throughput approach was 12.5% for designing novel PHA synthases. This enzyme is very challenging to engineer, as it is poorly characterized and requires dimerization for its function. The conditional variational autoencoder generated a diverse set of PHA synthase sequences by learning from native proteins, and we applied a series of tests to select the most promising candidates. We believe that the selection process may have contributed to the high success rate. However, our subsequent analysis revealed that the low frequency of unnatural amino acid pair combinations was a key feature of the active sequences, and we suggest using this criterion for future selections.

Materials and methods

Data collection and representation

Collection and processing of PHA synthase data. We used sequences from Uniprot [45] (uniprot version 2021_03) to create a data set for this study. These included all sequences containing class I, II, or III PHA synthase domains (InterPro families IPR010963, IPR011287, IPR010125 [46]), or N-terminal domain of a PHA synthase (InterPro family IPR010941 [46]). A class was assigned for each sequence based on InterPro family and a phylogenetic tree. ClustalW 2.1 [47] was used to create a multiple sequence alignment (MSA) of sequences belonging to InterPro [46] classes IPR010963, IPR011287 and IPR010125. The D3 JavaScript library [48] was used to create a phylogenetic tree of the MSA. A cluster with known class IV poly(R)-hydroxyalkanoate synthase sequences (Q8GI81 and Q9ZF92, Uniprot) [6] was assigned as class IV. The sequences assigned to class IV belonged to class III polyhydroxyalkanoate synthase family in InterPro (IPR010125). Thus, the sequences belonging to IPR010963 were assigned to class I, the sequences belonging to IPR011287 were assigned to class II, and the sequences belonging to IPR010125 were assigned to class III or IV based on the MSA and phylogenetic tree explained above. For sequences belonging to IPR010941 but not to IPR010963, IPR011287, or IPR010125, the class was assigned based on Basic Local Alignment Search Tool (BLAST). BLAST 2.12.0 [22] was run for each sequence against PhaC sequences belonging to IPR010963, IPR011287 or IPR010125. The class of the query sequence was assigned to be the same class as the PHA synthase sequence with lowest e-value. Furthermore, we augmented the data set with sequences with known beneficial mutations (S1 Table). Then, the sequences were classified based on their size. Very long sequences were removed and the rest of the sequences were divided to three different size categories (i.e., normal size, small, and fraction) based on the sequence length (S2 Table).

In addition to the PHA synthase sequences, 5000 lipase sequences were included in the dataset. Lipases were added to increase the size of the dataset as they have similar structures as the PHA synthases. Both lipases and PHA synthases contain α/β -hydrolase domains. Furthermore, PHA synthases contain lipase box sequence (GXSXG) with only one amino acid substitution compared to lipases. In PHA synthases, the active site serine is replaced with cysteine [23]. First, sequences containing lipase in their name were downloaded from Uniprot (v. 2021_3). Then BLAST 2.12.0 was run for these sequences against PhaC sequences belonging to IPR010963, IPR011287, or IPR010125, and 5000 lipase sequences with smallest e-value were included in the dataset.

Finally, the dataset was divided to training and testing datasets. The CD-HIT 4.6 [21] was used to cluster all sequences in each class (i.e., Class I-IV PHA synthases and lipases) separately. One cluster in each class was then selected to testing dataset. The data processing pipeline is clarified in S5 Fig. The number of sequences in different classes, the number of sequences in different size categories as well as the sizes of training and testing datasets are compiled in S3 Table. The dataset that was used to generate the TFRecord files and the TFRecord files are stored to Zenodo (<https://doi.org/10.5281/zenodo.17549219>).

Protein representation. We used a combination of features to represent the sequence and structure information of the enzymes for the generative model. To obtain structural features for our large dataset of enzymes, we used the software trRosetta [49]. It predicts the contact maps of each enzyme, which consist of a matrix of distance probabilities and three matrices of orientation probabilities for each pair of residues. Moreover, we encoded each protein sequence using seven physico-chemical amino acid properties, such as mass, side chain volume, or polarity (see Table 4), and the secondary structure of the protein predicted with STRIDE [50].

Building and use of a variational autoencoder (VAE)

Model architecture. Our model follows an autoencoder architecture, which consists of two parts: the encoder that compresses the input into a latent vector and the decoder that reconstructs the input from the compressed latent vector

Table 4. Features collected from each amino acid.

Property	Definition	Reference
Mass	Masses of neutral residues	[51]
Side chain volume	Mean volumes of residues inside proteins	[52]
Hydropathy	Hydrophobicity/hydrophilicity of the residues	[53]
Polarity	Average separation of charge in the residues	[54]
Polarizability	Possibility of the residues to become polarized temporally	[55]
pI	pH of the residues at the isoelectric point	[56]
Side chain composition	Atomic weight ratio of noncarbon elements inside side chain end groups or rings	[54]

<https://doi.org/10.1371/journal.pcbi.1014087.t004>

created by the encoder. However, our model differs from the standard autoencoders because it does not use the same input and output. Our encoder takes as input the structural and sequence-based features of enzymes, while the decoder outputs the one-hot representation of the amino acids in the enzyme sequence. An additional input to both the encoder and the decoder is the condition that represents each enzyme class (a 85-dimensional vector). The enzyme data for the model training comprised PHA synthases (4 classes), lipases (1 class), small PHA synthase sequences (4 classes), and partial PHA synthase sequences (4 classes) (S3 Table). We separated the small and partial sequences based on length and treated them as distinct classes in the condition vector.

We compared the convolutional [57] and recurrent (LSTM [58] and GRU [59]) model architectures and selected the LSTM-based model because it achieved the highest test accuracy. Autoencoder architectures consist of two parts: the encoder, which compresses the input given for the training, and the decoder, which reconstructs the input from compressed values produced by the encoder. In our model, the encoder input comprised the contact map representation of the protein (a 700x75 matrix), the protein amino acid features representation (a 700x7 matrix) (Table 4), the secondary structure prediction (a 700x3 matrix), and the condition vector (an 85-dimensional vector). Each of these inputs, except the condition, was passed through a multi-head attention block [60] with 18 heads, followed by a bidirectional LSTM layer. The outputs of these three blocks were then concatenated with the condition vector and fed into two dense layers to produce the mean and the logvar values corresponding to the input.

We sampled the values of the latent vector using the mean and the logvar values generated by the encoder. The decoder had two inputs: the latent vector sampled from the encoder (a 20-dimensional vector) and the same condition vector used in the encoder. Both inputs were concatenated and connected to a bidirectional LSTM layer. The output of the LSTM layer was flattened and sent to two consecutive dense layers that returned a one-hot representation of the amino acids in the protein sequence (a 700x21 matrix) (see Fig 1). Python implementation of the model is available online https://github.com/vttresearch/PHA_cVAE/.

Model training and evaluation. The model was trained using Adam optimizer with learning rate 0.001 and a batch size of 12. Softmax cross entropy with logits was used to calculate the reconstruction error in the loss function. To mitigate posterior collapse, the beta parameter, that controls the weight of KL divergence on the total loss, was incrementally increased from 0.01 to 1. Further increments in the beta parameter resulted in model collapse. First training was done for 203 epochs using beta parameter 0.01, then 214 epochs using beta parameter 0.1 and finally for 202 epochs using beta parameter 1. In each training stage training was continued sufficient time for the model not to improve anymore, but being sure that the model did not start overfitting. Convergence plots of loss, accuracy and KL divergence loss are presented in S6 Fig. Although validation accuracy did not improve during the last two training stages decrease in KL divergence loss, especially during the last training phase, facilitate novel protein generation by regularizing the latent space. Hyperparameter optimization was done using a grid search methodology. We tested 0.1, 0.01 and 0.001 values as the learning rate with the smallest learning rate being selected due to its better performance in terms of accuracy. Additionally,

we tested various dropout regularization values, but observed no significant performance differences. Furthermore, we evaluated different model parameters such as number of attention heads in the attention block, size of latent vector and amount of nodes in the dense layers of the encoder (S4 Table and S7 Fig). The batch size of 12 was limited by the memory of our GPU. The training took 18 days using NVIDIA Volta V100 GPU.

Model performance was evaluated by analyzing quality of the produced sequences. We compared the sequences generated with PHA_cVAE with sequences generated with a LLM and native class I PHA synthase sequences with 100 random mutations. As LLM we used the ProGen2-small [17] and followed their fine-tuning protocol [61]. The input data for the ProGen2-small was the same data as for our PHA_cVAE model, divided into 13 groups (12 categories for PHA synthases and lipases (S2 Table, S3 Table)). We fine-tuned the ProGen2-small for 5 epochs because the test error increased after the fifth epoch. We generated the amino acid sequences with the ProGen2-small using the label of the Class I PHA synthases followed by “M” (the first amino acid in all the proteins) as a seed.

We first evaluated the amino acid composition of the generated sequences by selecting native sequences from each PHA synthase class randomly (1000 from classes I, II and III and 300 from class IV) (InterPro families IPR010963, IPR011287, IPR010125) and counting the amount of each amino acid in these sequences. By dividing the amount of each amino acid with the length of the sequence we defined the percentage range of each amino acid in the native sequences. We then generated 1000 new PHA synthase sequences from each class with PHA_cVAE and with ProGen2-small and analyzed how many of them had all the amino acids within the calculated range.

Next, we predicted structures with Boltz-1 [18] for 750 sequences generated with each generative model (ProGen2-small and PHA_cVAE) and 495 randomly mutated native PHA synthase sequences that had been filtered according to the amino acid composition. We first selected 750 sequences randomly from 10322 sequences that had been generated with PHA_cVAE and already were filtered according to the amino acid composition (Section “Sequence generation and selection of the 16 sequences”). Next, we extracted 750 sequences that had similar amino acid composition as native PHA synthases from 1000 sequences generated with ProGen2-small. In addition, we randomly selected 750 native class I PHA synthase sequences from the training and test dataset, generated randomly 100 mutations to each and filtered the sequences according to the amino acid composition leaving us 495 randomly mutated native PHA synthase sequences. We then calculated the average pLDDT score [19] for the structures predicted with Boltz-1. In addition, we aligned the structures with the PHA synthase from *Chromobacterium* sp. USM2 using TM-align [20], as we didn't expect a significant structural divergence to this protein, and calculated the average TM-score [20]. Using the TM-align alignments, we also checked the proportion of generated PHA synthases that didn't contain the three catalytic amino acids in the correct position.

Finally, we evaluated if PHA synthase classes of the generated sequences were correct and whether PHA_cVAE was confident about terminating the sequence. Evaluating the classes of the generated sequences were done by generating 100 sequences in each class with PHA_cVAE, running BLAST 2.12.0 against native PhaC sequences (belonging to InterPro families IPR010963, IPR011287 or IPR010125) and evaluating if the hit with the lowest e-value was from the desired class. When generating sequences for class evaluation only sequences with similar amino acid composition as native sequences were saved. Class evaluation was also made similarly for 750 ProGen2-small generated sequences. Analysis of sequence termination was performed by generating 1000 sequences in each class with PHA_cVAE and evaluating whether the sequences contained termination indicators elsewhere than in the end of the sequence.

Sequence generation and selection of the 16 sequences. We generated 10322 different class I PHA synthase sequences using the selected cVAE generative model and by randomizing the values in the latent space. More than 10 000 sequences were generated to have enough variance in the generated sequences before starting filtering process. When generating sequences, only the sequences having amino acid proportions within the same range as in the native sequences (calculated similarly as in Materials and methods, Model training and evaluation) were saved. Next we removed duplicated sequences and filtered the generated PHA synthases based on the presence/absence of the three

catalytic amino acids in the right position in active site leaving 9241 enzymes. With the remaining sequences, we used CD-HIT [21] with 95% identity threshold to cluster them resulting on 715 clusters. When generating sequences PHA_cVAE gave the probability of each amino acid for each position in the sequence. The amino acid with highest probability was then selected to the generated sequence. To select sequences from the clusters we analyzed for each position in the sequence, how many amino acids had similar probability as the selected amino acid in each position of the sequence (maximum 10% smaller probability than the amino acid selected to the sequence). Then from each cluster we selected the sequence with least of amino acids having similar probabilities than the amino acid with highest probability as we considered that the model was less confident in those positions.

Next the filtering continued by running AlphaFold2 [19] for each of the 715 sequences and aligning the obtained structures with the crystal structure of PHA synthase from *Chromobacterium* sp. USM2 (PDBe: 5xav [23]) using TM-align [20]. These aligned structures were then used to analyze if the created new PHA synthase sequences had changes in conserved residues (other than the catalytic triad) obtained from Chek *et al.* [23] (amino acids corresponding to 197L, 200Y, 211P, 213L, 220N, 223Y, 226D, 232S, 249W, 289G, 293G, 294G, 323D, 365R, 392W, 395D, 415N, 431D, 448H, 476G, 489K in *Chromobacterium* USM2). In addition, we studied if amino acids in the dimer interface corresponding to amino acids at positions 332, 333, 369, 371, 386, 387, 390, and 451 in *Chromobacterium* sp. USM2 [23] were hydrophobic. However, the alignment of our structures with 5xav was not good for positions between residues 371 and 386 in PHA synthase of *Chromobacterium* sp. USM2. We expect this to be due to the break in the crystal structure at positions 372–384. Thus, we aligned the generated 715 sequences also with the AlphaFold structure of the PHA synthase from *Chromobacterium* sp. USM2 (AlphaFoldDB: E1APK1) and used these alignments to check the hydrophobicity of the amino acids corresponding to amino acids at positions 372–384 of E1APK1.

Next we analyzed the amount of α -helices in the N-terminal part of the PHA_cVAE generated PHA synthase sequences. We aligned AlphaFold predicted structures with PHA synthase from *Cupriavidus necator* (AlphaFoldDB: P23608) using TM-align, analyzed the start of the N-terminal from the aligned structures and predicted the secondary structure from the AlphaFold predicted structures using STRIDE [50].

Furthermore, we checked the lengths of the generated sequences and sequence termination indicators inside the sequences. For sequences containing a termination indicator inside the sequence, the termination indicator was changed to the amino acid with second highest probability. With the information explained above, we then selected 42 sequences for further analysis. For these 42 sequences, we calculated the tunnels to and from the active site with CAVER 3.0.3 Pymol plugin [25] and selected 16 to be tested in wet lab using all the information gathered above. Finally, we aligned the sequences containing a termination indicator inside the sequence with their closest native PHA synthase and analyzed if the native PHA synthase was longer or shorter. Closest native sequence of PhaC_{VAE10} finished at the position where PhaC_{VAE10} had first termination indicator. Thus, we decided to cut this sequence at that position. Rest of the native sequences were longer than the PHA_cVAE generated sequences and therefore, we kept these sequences as they were.

To analyze if the selection process led us to select sequences uniformly from all the produced sequences, we generated a phylogenetic tree of the PHA_cVAE generated sequences (Fig 2). First, we clustered all class I PHA synthase sequences generated with the model using CD-HIT [21] using 50% identity threshold. This resulted in one representative sequence for each cluster. Next, all 16 sequences selected for *in vivo* experiments were added to the set of sequences. BLAST 2.15.0 [22] was then run against native PHA synthase sequences (belonging to InterPro families IPR010963, IPR011287 or IPR010125) and the hit with lowest percentage identity was selected. A phylogenetic tree of the PHA_cVAE generated sequences was then generated using ClustalW 2.1 [47] and visualised using Geneious 10.2.6 (<https://www.geneious.com>). The closest native sequences were marked in the branch labels.

Latent space visualization

Latent vectors of the sequences selected for wet lab validation were visualized together with 1000 randomly selected PHA_cVAE generated sequences and 1674 randomly selected native PHA synthase sequences. Each sequence was

presented with the same protein representation that was used during model training and the inputs were passed through the encoder to obtain the latent vectors. The latent vectors of 20 dimensions were then visualized in 2D using UMAP. Jupyter Notebook containing the visualization can be found in [S1 File](#).

In vivo activity measurement of novel PHA syntases

Strain engineering. Strains and enzymes used in this study are presented in [Table 3](#). The used plasmids are listed in [S5 Table](#). *S. cerevisiae* CEN.PK113-7D was kindly provided by Dr. P. Kötter (Institut für Mikrobiologie, J.W. Goethe Universität Frankfurt, Germany). Pathway to produce 3-HB-CoA (*phaA* and *phaB1*) was cloned to three different EasyClone vectors pCfB3034, pCfB2903, and pCfB3039 by amplifying the precursor pathway (*pTEF1-phaA-tENO1-pTDH3-phaB*) from plasmid pPHB_{template_1} (B9660) [62] and cloning the product into the corresponding EasyClone plasmids in front of *tCYC1* terminator using Gibson assembly (E2611S, New England BioLabs). The generated plasmids were digested with NotI and transformed to CEN.PK113-7D (H3887) to generate strain PhaA-PhaB1 (H6135). Gene coding for *phaC1_{Cn}* (*pPGK1-phaC1_{Cn}*) was amplified with PCR from pPHB_{template_1} (B9660) and cloned to EasyClone vector pCfB3035 in front of *tCYC1* terminator with Gibson assembly. Genes coding for VAE generated novel PHA synthases (PHAC_{VAE1}-PHAC_{VAE16}) and genes coding for PhaC_{LS} and PhaC_{Jl} were codon optimized for *S. cerevisiae* and ordered with PGK1 promoters from GenScript in pCfB3035 EasyClone vectors. The plasmids were digested with NotI and transformed to parent strain PhaA-PhaB1 (H6135) to generate strains PHA_VAE1-VAE16. Strain PHB_ctr2 was built by amplifying PHB pathway (*pTEF1-phaA-tENO1-pTDH3-phaB1-tSSA1-pPGK1-phaC1_{Cn}-tCYC*) with PCR from the pPHB_{template_2} (B11787) and cloning the product to EasyClone vectors pCfB3035, pCfB2904, and pCfB2909 using Gibson assembly. These vectors were then digested with NotI and transformed to CEN.PK113-7D (H3887). Lithium acetate (LiAc)/ SS carrier DNA/ PEG method [63] and CRISPR/Cas9 protocol of the EasyClone kit [27] were used in all transformations. Correct integrations were confirmed with PCR using oligos of EasyClone kit and gene specific oligos as well as Sanger sequencing (Microsynth Seqlab GmbH).

Nile red analysis. The strains were grown for 16 h in 3 ml of synthetic complete media with 20 g/l of glucose in a 24 well plate at 770 rpm shaking and 30 °C. Cultivation was started by inoculating the media with the strains from YPD plates. At the end of the cultivation each cell culture was diluted with distilled water to obtain OD₆₀₀ 2. Then 100 µl of each of the diluted samples was transferred to a black 96-well plate and mixed with 20 µl of Nile red dissolved in DMSO, so that the final concentrations of Nile red and DMSO were 5 mg/l and 17% v/v, respectively. Fluorescence was measured after 10 min of incubation at RT with Varioskan Flash (Thermo scientific) using 550 nm excitation and 610 nm emission wavelengths. Each culture was measured in four technical replicates. The protocol is available with DOI: <https://dx.doi.org/10.17504/protocols.io.5jyl8xwxrv2w/v1>.

Shake flask cultivation. All the strains were grown in synthetic complete media with 20 g/l of glucose at 30 °C and 220 rpm shaking. Precultures of 10 ml in 50 ml Erlenmeyer flasks were grown overnight. Subsequently 50 ml cultures in 250 ml Erlenmeyer flasks were started from OD₆₀₀ 0.2 and continued for 72 hours. Samples were taken at 24 h and 72 h to monitor the population growth, extracellular metabolite formation, and glucose utilization. At the end of the 72 h cultivation, cells were harvested by centrifuging them for 6 min at 4000 rpm and washing once with distilled water. However, the washing step was not done for strains PHA_VAE3_{A-C}, PHA_VAE4_{A-C}, PHA_VAE5_{A-B}, PHA_VAE6_A, PHA_VAE7_A and PHA_VAE9_{A-B}. Here A, B, and C refers to three different biological replicates of each strain. The cell pellets were stored in -20°C until GC-MS analysis. The population growth was analyzed by measuring the optical density (OD₆₀₀).

PHB quantitation with GC-MS. The amount of accumulated PHB in CDW was analyzed with gas chromatography mass spectrometry (GC-MS) similarly to Ylinen *et al.* [30] based on method described by Braunegg *et al.* [64]. Cell samples from the shake flask cultivations were frozen at -80 °C and lyophilized overnight in Christ Alpha 2–4 LSCBasic device. Ten milligrams of each dried sample was subjected to methanolysis for 140

min in 100 °C water bath in a solution containing 1 ml chloroform, 810 µl methanol, 150 µl sulfuric acid, and 50 µl 3-hydroxybutyric acid 1,3-¹³C₂ (Sigma-Aldrich) as an internal standard. Samples were cooled down to room temperature and 1 ml of distilled water was added to remove water-soluble particles. Chloroform phase was then analyzed using gas chromatography equipment (7890, Agilent) with HP-FFAP column (19091F-102 Agilent). Two replicates were analyzed of each strain. A 3-hydroxybutyric acid standard was analyzed equally as the samples. Quantitative results for PHB content were corrected using a recovery value of 91.5% and by considering the molecular weight difference (18 amu) between the monomer unit in the polymer and the free acid used in the calibration curve. The recovery value was obtained by assessing average recovery % of commercial PHB polymer at two concentration levels 500 µg and 2000 µg corresponding to 5% and 20% PHB content in a 10 mg biological sample, respectively. Recovery values for 500 µg and 2000 µg samples were 91.6% and 91.4% with 1.5% and 1.2% relative standard deviations (RSDs), respectively.

Sequence analysis of the generated novel PHA synthases

Generated PHA synthase sequences were BLASTed against all the PHA synthase sequences collected from Uniprot [45] to identify the most similar native PHA synthase for each generated novel PHA synthase and percentage of identical positions were collected to Table 2. Then generated PHA synthases were structurally aligned (structures were predicted with AlphaFold2 [19]) to their most similar native PHA synthase enzymes using TM-align [20] and the positions of amino acid differences as well as amino acids in the generated enzyme at these positions were recorded. We then analyzed how rare it was to have amino acids found in the novel PHA synthases at these positions. For each position we collected all native class I PHA synthase enzymes with the same amino acid as the novel PHA synthase in that position. If the number of collected native class I PHA synthases were less than one hundred we considered the amino acid at that position as a rare amino acid and if no native class I PHA synthases were found we considered the change as a novel change. To do this analysis all 16 generated novel PHA synthases and all native class I PHA synthases were structurally aligned with class I PHA synthase from *Cupriavidus necator* using TM-align [20]. This alignment was used to align all the PHA synthases with each other (i.e., to know which amino acid position in a native PHA synthase corresponds to the amino acid position under inspection in the novel PHA synthase).

Next we analyzed pairwise combinations of amino acid differences. For each generated PHA synthase we collected all pairs of positions where different amino acids were found with respect to the most similar native class I PHA synthase. We then looked for these pairs in native class I PHA synthases and counted the cases where there was no native enzyme with the same two amino acids as the novel PHA synthase under inspection (see Table 2).

To generate Fig 4 and S3 Fig we visualized and aligned the structures obtained with AlphaFold with PyMOL 3.2.0a and marked amino acid differences.

Supporting information

S1 Fig. Distribution plots for pLDDT scores and TM-scores.

(PDF)

S2 Fig. Growth of the strains during shake flask cultivations.

(PDF)

S3 Fig. Structural alignment of PhaC_{VAE2} and PhaC_{VAE6} with their closest native PHA synthases.

(PDF)

S4 Fig. Analysis of CAVER tunnels.

(PDF)

S5 Fig. Processing training data.

(PDF)

S6 Fig. Convergence plots of loss, accuracy, and KL divergence loss.

(PDF)

S7 Fig. Performance comparison of evaluated models.

(PDF)

S1 Table. Mutated PHA synthases added to the training dataset.

(PDF)

S2 Table. Number of amino acids (aa) in the different size categories.

(PDF)

S3 Table. Number of sequences in different size categories, training dataset, and testing dataset.

(PDF)

S4 Table. Different parameters used when selecting the model.

(PDF)

S5 Table. Plasmids used in this study.

(PDF)

S1 File. Latent space visualization.

(IPYNB)

S1 Data. Data for latent space visualization.

(CSV)

Acknowledgments

We acknowledge the invaluable support of Kaisa Peltonen for storing strains to culture collection and Matti Hölttä for GC-MS analysis. In addition, we acknowledge Samuli Ollila, Heli Nygren and Gopal Peddinti for their comments on the manuscript. Their assistance has been crucial to our efforts. Furthermore, this work was carried out under the Academy of Finland Center of Excellence Program (2022–2029) in Life-Inspired Hybrid Materials (LIBER), project number (346106).

Author contributions

Conceptualization: Tuula Marjaana Tenkanen, Anna Ylinen, Paula Jouhten, Merja Penttilä, Sandra Castillo.

Formal analysis: Tuula Marjaana Tenkanen, Sandra Castillo.

Funding acquisition: Merja Penttilä.

Investigation: Tuula Marjaana Tenkanen, Sandra Castillo.

Methodology: Sandra Castillo.

Software: Tuula Marjaana Tenkanen, Sandra Castillo.

Supervision: Anna Ylinen, Merja Penttilä, Sandra Castillo.

Visualization: Tuula Marjaana Tenkanen, Sandra Castillo.

Writing – original draft: Tuula Marjaana Tenkanen, Anna Ylinen, Sandra Castillo.

Writing – review & editing: Tuula Marjaana Tenkanen, Anna Ylinen, Paula Jouhten, Sandra Castillo.

References

1. Neoh Z, Chek F, Tan T, Linares-Pastén JA, Nandakumar A, Hakoshima T. Polyhydroxyalkanoate synthase (PhaC): the key enzyme for biopolyester synthesis. *Curr Res Biotechnol*. 2022;4:87–101. <https://doi.org/10.1016/j.crbiot.2022.01.002>
2. Choi SY, Cho IJ, Lee Y, Kim YJ, Kim KJ, Lee SY. Microbial polyhydroxyalkanoates and nonnatural polyesters. *Adv Mater*. 2020;32(35):1907138. <https://doi.org/10.1002/adma.201907138>
3. Chek MF, Kim SY, Mori T, Matsumoto K, Sato S, Hakoshima T. Structures of polyhydroxyalkanoate synthase PhaC from *Aeromonas caviae*, producing biodegradable plastics. *Angew Chem Int Ed*. 2025;64(26):e202504626. <https://doi.org/https://doi.org/10.1002/anie.202504626>
4. Assefa NG, Hansen H, Altermark B. A unique class I polyhydroxyalkanoate synthase (PhaC) from *Brevundimonas* sp. KH11J01 exists as a functional trimer: a comparative study with PhaC from *Cupriavidus necator* H16. *N Biotechnol*. 2022;70:57–66. <https://doi.org/10.1016/j.nbt.2022.05.003> PMID: 35533829
5. Li P, Chakraborty S, Stubbe J. Detection of covalent and noncovalent intermediates in the polymerization reaction catalyzed by a C149S class III polyhydroxybutyrate synthase. *Biochemistry*. 2009;48(39):9202–11. <https://doi.org/10.1021/bi901329b> PMID: 19711985
6. Rehm BHA. Polyester synthases: natural catalysts for plastics. *Biochem J*. 2003;376(Pt 1):15–33. <https://doi.org/10.1042/BJ20031254> PMID: 12954080
7. Khersonsky O, Lipsh R, Avizemer Z, Ashani Y, Goldsmith M, Leader H. Automated design of efficient and functionally diverse enzyme repertoires. *Mol Cell*. 2018;72(1):178–186.e5. <https://doi.org/10.1016/j.molcel.2018.08.033>
8. Lipsh-Sokolik R, Khersonsky O, Schröder SP, de Boer C, Hoch S-Y, Davies GJ, et al. Combinatorial assembly and design of enzymes. *Science*. 2023;379(6628):195–201. <https://doi.org/10.1126/science.ade9434> PMID: 36634164
9. Listov D, Goverde CA, Correia BE, Fleishman SJ. Opportunities and challenges in design and optimization of protein function. *Nat Rev Mol Cell Biol*. 2024;25(8):639–53. <https://doi.org/10.1038/s41580-024-00718-y> PMID: 38565617
10. Sevgen E, Moller J, Lange A, Parker J, Quigley S, Mayer J, et al. ProT-VAE: Protein Transformer Variational AutoEncoder for functional protein design. *Proc Natl Acad Sci U S A*. 2025;122(41):e2408737122. <https://doi.org/10.1073/pnas.2408737122> PMID: 41052325
11. Hawkins-Hooker A, Depardieu F, Baur S, Couairon G, Chen A, Bikard D. Generating functional protein variants with variational autoencoders. *PLoS Comput Biol*. 2021;17(2):e1008736. <https://doi.org/10.1371/journal.pcbi.1008736> PMID: 33635868
12. Greener JG, Moffat L, Jones DT. Design of metalloproteins and novel protein folds using variational autoencoders. *Sci Rep*. 2018;8(1):16189. <https://doi.org/10.1038/s41598-018-34533-1> PMID: 30385875
13. Ylinen A, Salusjärvi L, Toivari M, Penttilä M. Control of D-lactic acid content in P(LA-3HB) copolymer in the yeast *Saccharomyces cerevisiae* using a synthetic gene expression system. *Metab Eng Commun*. 2022;14:e00199. <https://doi.org/10.1016/j.mec.2022.e00199> PMID: 35571351
14. Joyce JM. Kullback-Leibler divergence. In: Lovric M, editor. *International Encyclopedia of Statistical Science*. Berlin, Heidelberg: Springer. 2011. p. 720–2.
15. Higgins I, Matthey L, Pal A, Burgess C, Glorot X, Botvinick M. β -VAE: Learning basic visual concepts with a constrained variational framework. 2017. <https://openreview.net/forum?id=Sy2fzU9gI>
16. Ichikawa Y, Hukushima K. Dataset Size Dependence of Rate-Distortion Curve and Threshold of Posterior Collapse in Linear VAE. 2023. <https://arxiv.org/abs/2309.07663>
17. Nijkamp E, Ruffolo JA, Weinstein EN, Naik N, Madani A. ProGen2: Exploring the boundaries of protein language models. *Cell Syst*. 2023;14(11):968–978.e3. <https://doi.org/10.1016/j.cels.2023.10.002> PMID: 37909046
18. Wohlwend J, Corso G, Passaro S, Reveiz M, Leidal K, Swiderski W. Boltz-1 Democratizing Biomolecular Interaction Modeling. 2024. <https://www.biorxiv.org/content/10.1101/2024.11.19.624167v1>
19. Jumper J, Evans R, Pritzel A, Green T, Figurnov M, Ronneberger O, et al. Highly accurate protein structure prediction with AlphaFold. *Nature*. 2021;596(7873):583–9. <https://doi.org/10.1038/s41586-021-03819-2> PMID: 34265844
20. Zhang Y, Skolnick J. TM-align: a protein structure alignment algorithm based on the TM-score. *Nucleic Acids Res*. 2005;33(7):2302–9. <https://doi.org/10.1093/nar/gki524> PMID: 15849316
21. Li W, Godzik A. Cd-hit: a fast program for clustering and comparing large sets of protein or nucleotide sequences. *Bioinformatics*. 2006;22(13):1658–9. <https://doi.org/10.1093/bioinformatics/btl158> PMID: 16731699
22. Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, et al. BLAST+: architecture and applications. *BMC Bioinformatics*. 2009;10:421. <https://doi.org/10.1186/1471-2105-10-421> PMID: 20003500
23. Chek MF, Kim S-Y, Mori T, Arsad H, Samian MR, Sudesh K, et al. Structure of polyhydroxyalkanoate (PHA) synthase PhaC from *Chromobacterium* sp. USM2, producing biodegradable plastics. *Sci Rep*. 2017;7(1):5312. <https://doi.org/10.1038/s41598-017-05509-4> PMID: 28706283
24. Kim Y-J, Choi SY, Kim J, Jin KS, Lee SY, Kim K-J. Structure and function of the N-terminal domain of *Ralstonia eutropha* polyhydroxyalkanoate synthase, and the proposed structure and mechanisms of the whole enzyme. *Biotechnol J*. 2017;12(1):10.1002/biot.201600649. <https://doi.org/10.1002/biot.201600649> PMID: 27808475
25. Chovancova E, Pavelka A, Benes P, Strnad O, Brezovsky J, Kozlikova B, et al. CAVER 3.0: a tool for the analysis of transport pathways in dynamic protein structures. *PLoS Comput Biol*. 2012;8(10):e1002708. <https://doi.org/10.1371/journal.pcbi.1002708> PMID: 23093919

26. Asperti A. Variance loss in variational autoencoders. *Lecture Notes in Computer Science*. Springer; 2020. p. 297–308. https://doi.org/10.1007/978-3-030-64583-0_28
27. Jessop-Fabre MM, Jakočiūnas T, Stovicek V, Dai Z, Jensen MK, Keasling JD, et al. EasyClone-MarkerFree: A vector toolkit for marker-less integration of genes into *Saccharomyces cerevisiae* via CRISPR-Cas9. *Biotechnol J*. 2016;11(8):1110–7. <https://doi.org/10.1002/biot.201600147> PMID: [27166612](https://pubmed.ncbi.nlm.nih.gov/27166612/)
28. Zuriani R, Vigneswari S, Azizan MNM, Majid MIA, Amirul AA. A high throughput Nile red fluorescence method for rapid quantification of intracellular bacterial polyhydroxyalkanoates. *Biotechnol Bioproc E*. 2013;18(3):472–8. <https://doi.org/10.1007/s12257-012-0607-z>
29. Kacmar J, Carlson R, Balogh SJ, Srienc F. Staining and quantification of poly-3-hydroxybutyrate in *Saccharomyces cerevisiae* and *Cupriavidus necator* cell populations using automated flow cytometry. *Cytometry A*. 2006;69(1):27–35. <https://doi.org/10.1002/cyto.a.20197> PMID: [16342115](https://pubmed.ncbi.nlm.nih.gov/16342115/)
30. Ylinen A, de Ruijter JC, Jouhten P, Penttilä M. PHB production from cellobiose with *Saccharomyces cerevisiae*. *Microb Cell Fact*. 2022;21(1):124. <https://doi.org/10.1186/s12934-022-01845-x> PMID: [35729556](https://pubmed.ncbi.nlm.nih.gov/35729556/)
31. Sandström AG, Muñoz de Las Heras A, Portugal-Nunes D, Gorwa-Grauslund MF. Engineering of *Saccharomyces cerevisiae* for the production of poly-3-d-hydroxybutyrate from xylose. *AMB Express*. 2015;5:14. <https://doi.org/10.1186/s13568-015-0100-0> PMID: [25852991](https://pubmed.ncbi.nlm.nih.gov/25852991/)
32. Breuer U, Terentiev Y, Kunze G, Babel W. Yeasts as producers of polyhydroxyalkanoates: genetic engineering of *Saccharomyces cerevisiae*. *Macromol Biosci*. 2002;2(8):380–6. [https://doi.org/10.1002/1616-5195\(200211\)2:8<380::AID-MABI380>3.0.CO;2-X](https://doi.org/10.1002/1616-5195(200211)2:8<380::AID-MABI380>3.0.CO;2-X)
33. Carlson R, Srienc F. Effects of recombinant precursor pathway variations on poly[(R)-3-hydroxybutyrate] synthesis in *Saccharomyces cerevisiae*. *J Biotechnol*. 2006;124(3):561–73. <https://doi.org/10.1016/j.jbiotec.2006.01.035> PMID: [16530287](https://pubmed.ncbi.nlm.nih.gov/16530287/)
34. de Las Heras AM, Portugal-Nunes DJ, Rizza N, Sandström AG, Gorwa-Grauslund MF. Anaerobic poly-3-D-hydroxybutyrate production from xylose in recombinant *Saccharomyces cerevisiae* using a NADH-dependent acetoacetyl-CoA reductase. *Microb Cell Fact*. 2016;15(1):197. <https://doi.org/10.1186/s12934-016-0598-0> PMID: [27863495](https://pubmed.ncbi.nlm.nih.gov/27863495/)
35. Portugal-Nunes DJ, Pawar SS, Lidén G, Gorwa-Grauslund MF. Effect of nitrogen availability on the poly-3-D-hydroxybutyrate accumulation by engineered *Saccharomyces cerevisiae*. *AMB Express*. 2017;7(1):35. <https://doi.org/10.1186/s13568-017-0335-z> PMID: [28176283](https://pubmed.ncbi.nlm.nih.gov/28176283/)
36. Christensen M, Chiciudean I, Jablonski P, Tanase A-M, Shapaval V, Hansen H. Towards high-throughput screening (HTS) of polyhydroxyalkanoate (PHA) production via Fourier transform infrared (FTIR) spectroscopy of *Halomonas* sp. R5-57 and *Pseudomonas* sp. MR4-99. *PLoS One*. 2023;18(3):e0282623. <https://doi.org/10.1371/journal.pone.0282623> PMID: [36888636](https://pubmed.ncbi.nlm.nih.gov/36888636/)
37. Yang TH, Jung YK, Kang HO, Kim TW, Park SJ, Lee SY. Tailor-made type II *Pseudomonas* PHA synthases and their use for the biosynthesis of polylactic acid and its copolymer in recombinant *Escherichia coli*. *Appl Microbiol Biotechnol*. 2011;90(2):603–14. <https://doi.org/10.1007/s00253-010-3077-2> PMID: [21221571](https://pubmed.ncbi.nlm.nih.gov/21221571/)
38. Pederson EN, McChalicher CWJ, Srienc F. Bacterial synthesis of PHA block copolymers. *Biomacromolecules*. 2006;7(6):1904–11. <https://doi.org/10.1021/bm0510101> PMID: [16768413](https://pubmed.ncbi.nlm.nih.gov/16768413/)
39. Wittenborn EC, Jost M, Wei Y, Stubbe J, Drennan CL. Structure of the Catalytic Domain of the Class I Polyhydroxybutyrate Synthase from *Cupriavidus necator*. *J Biol Chem*. 2016;291(48):25264–77. <https://doi.org/10.1074/jbc.M116.756833> PMID: [27742839](https://pubmed.ncbi.nlm.nih.gov/27742839/)
40. Kim J, Kim Y-J, Choi SY, Lee SY, Kim K-J. Crystal structure of *Ralstonia eutropha* polyhydroxyalkanoate synthase C-terminal domain and reaction mechanisms. *Biotechnol J*. 2017;12(1):10.1002/biot.201600648. <https://doi.org/10.1002/biot.201600648> PMID: [27808482](https://pubmed.ncbi.nlm.nih.gov/27808482/)
41. Du Z, Su H, Wang W, Ye L, Wei H, Peng Z, et al. The trRosetta server for fast and accurate protein structure prediction. *Nat Protoc*. 2021;16(12):5634–51. <https://doi.org/10.1038/s41596-021-00628-9> PMID: [34759384](https://pubmed.ncbi.nlm.nih.gov/34759384/)
42. Baek M, DiMaio F, Anishchenko I, Dauparas J, Ovchinnikov S, Lee GR, et al. Accurate prediction of protein structures and interactions using a three-track neural network. *Science*. 2021;373(6557):871–6. <https://doi.org/10.1126/science.abj8754> PMID: [34282049](https://pubmed.ncbi.nlm.nih.gov/34282049/)
43. Bowman SR, Vilnis L, Vinyals O, Dai AM, Jozefowicz R, Bengio S. Generating sentences from a continuous space. In: 2016. <http://arxiv.org/abs/1511.06349>
44. Austin J, Johnson DD, Ho J, Tarlow D, Berg Rvd. Structured denoising diffusion models in discrete state-spaces. 2023. <http://arxiv.org/abs/2107.03006>
45. The UniProt Consortium. UniProt: the universal protein knowledgebase in 2023. *Nucleic Acids Res*. 2023;51(D1):D523–31. <https://doi.org/10.1093/nar/gkac1052>
46. Paysan-Lafosse T, Blum M, Chuguransky S, Grego T, Pinto BL, Salazar G. InterPro in 2022. *Nucleic Acids Research*. 2022;51(D1):D418–27. <https://doi.org/10.1093/nar/gkac993>
47. Larkin MA, Blackshields G, Brown NP, Chenna R, McGettigan PA, McWilliam H. Clustal W and Clustal X version 2.0. *Bioinformatics*. 2007;23(21):2947–8. <https://doi.org/10.1093/bioinformatics/btm404>
48. Shank SD, Weaver S, Kosakovsky Pond SL. phylotree.js - a JavaScript library for application development and interactive data visualization in phylogenetics. *BMC Bioinformatics*. 2018;19(1):276. <https://doi.org/10.1186/s12859-018-2283-2> PMID: [30045713](https://pubmed.ncbi.nlm.nih.gov/30045713/)
49. Yang J, Anishchenko I, Park H, Peng Z, Ovchinnikov S, Baker D. Improved protein structure prediction using predicted interresidue orientations. *Proc Natl Acad Sci U S A*. 2020;117(3):1496–503. <https://doi.org/10.1073/pnas.1914677117> PMID: [31896580](https://pubmed.ncbi.nlm.nih.gov/31896580/)
50. Frishman D, Argos P. Knowledge-based protein secondary structure assignment. *Proteins*. 1995;23(4):566–79. <https://doi.org/10.1002/prot.340230412> PMID: [8749853](https://pubmed.ncbi.nlm.nih.gov/8749853/)

51. Voet D, Voet JG, Pratt CW. Voet's Principles of Biochemistry. Wiley. 2018.
52. Harpaz Y, Gerstein M, Chothia C. Volume changes on protein folding. *Structure*. 1994;2(7):641–9. [https://doi.org/https://doi.org/10.1016/S0969-2126\(00\)00065-4](https://doi.org/https://doi.org/10.1016/S0969-2126(00)00065-4)
53. Kyte J, Doolittle RF. A simple method for displaying the hydropathic character of a protein. *J Mol Biol*. 1982;157(1):105–32. [https://doi.org/10.1016/0022-2836\(82\)90515-0](https://doi.org/10.1016/0022-2836(82)90515-0) PMID: [7108955](https://pubmed.ncbi.nlm.nih.gov/7108955/)
54. Grantham R. Amino acid difference formula to help explain protein evolution. *Science*. 1974;185(4154):862–4. <https://doi.org/10.1126/science.185.4154.862> PMID: [4843792](https://pubmed.ncbi.nlm.nih.gov/4843792/)
55. Charton M, Charton BI. The structural dependence of amino acid hydrophobicity parameters. *J Theor Biol*. 1982;99(4):629–44. [https://doi.org/10.1016/0022-5193\(82\)90191-6](https://doi.org/10.1016/0022-5193(82)90191-6) PMID: [7183857](https://pubmed.ncbi.nlm.nih.gov/7183857/)
56. Brown TE, LeMay HE, Bursten BE, Murphy C, Woodward P, Stoltzfus ME. *Chemistry: The Central Science*. Pearson. 2019.
57. LeCun Y, Bengio Y, Hinton G. Deep learning. *Nature*. 2015;521(7553):436–44. <https://doi.org/10.1038/nature14539> PMID: [26017442](https://pubmed.ncbi.nlm.nih.gov/26017442/)
58. Hochreiter S, Schmidhuber J. Long short-term memory. *Neural Comput*. 1997;9(8):1735–80. <https://doi.org/10.1162/neco.1997.9.8.1735> PMID: [9377276](https://pubmed.ncbi.nlm.nih.gov/9377276/)
59. Cho K, van Merriënboer B, Gulcehre C, Bahdanau D, Bougares F, Schwenk H. Learning phrase representations using RNN encoder-decoder for statistical machine translation. In: 2014. <http://arxiv.org/abs/1406.1078>
60. Cordonnier JB, Loukas A, Jaggi M. Multi-head attention: collaborate instead of concatenate. 2021. <http://arxiv.org/abs/2006.16362>
61. Nijkamp E. enijkamp/progen2. <https://github.com/enijkamp/progen2>
62. Ylinen A, Maaheimo H, Anghelescu-Hakala A, Penttilä M, Salusjärvi L, Toivari M. Production of D-lactic acid containing polyhydroxyalkanoate polymers in yeast *Saccharomyces cerevisiae*. *J Ind Microbiol Biotechnol*. 2021;48(5–6):kuab028. <https://doi.org/10.1093/jimb/kuab028> PMID: [33899921](https://pubmed.ncbi.nlm.nih.gov/33899921/)
63. Gietz RD, Schiestl RH. High-efficiency yeast transformation using the LiAc/SS carrier DNA/PEG method. *Nat Protoc*. 2007;2(1):31–4. <https://doi.org/10.1038/nprot.2007.13> PMID: [17401334](https://pubmed.ncbi.nlm.nih.gov/17401334/)
64. Braunegg G, Sonnleitner B, Lafferty RM. A rapid gas chromatographic method for the determination of poly- β -hydroxybutyric acid in microbial biomass. *Eur J Appl Microbiol Biotechnol*. 1978;6(1):29–37. <https://doi.org/10.1007/BF00500854>