

# Absolute proteome quantification in the gas-fermenting acetogen *Clostridium autoethanogenum*

Kaspar Valgepea<sup>1,2</sup>, Gert Talbo<sup>1,3</sup>, Nobuaki Takemori<sup>4</sup>, Ayako Takemori<sup>4</sup>, Christina Ludwig<sup>5</sup>, Alexander P. Mueller<sup>6</sup>, Ryan Tappel<sup>6</sup>, Michael Köpke<sup>6</sup>, Séan Dennis Simpson<sup>6</sup>, Lars Keld Nielsen<sup>1,3,7</sup> and Esteban Marcellin<sup>1,3\*</sup>

<sup>1</sup>Australian Institute for Bioengineering and Nanotechnology (AIBN), The University of Queensland, 4072 St. Lucia, Australia

<sup>2</sup>ERA Chair in Gas Fermentation Technologies, Institute of Technology, University of Tartu, 50411 Tartu, Estonia

<sup>3</sup>Queensland Node of Metabolomics Australia, AIBN, The University of Queensland, 4072 St. Lucia, Australia

<sup>4</sup>Institute for Promotion of Science and Technology, Ehime University, 791-0295 Ehime, Japan

<sup>5</sup>Bavarian Center for Biomolecular Mass Spectrometry (BayBioMS), Technical University of Munich, 85354 Freising, Germany

<sup>6</sup>LanzaTech Inc., 60077 Skokie, USA

<sup>7</sup>The Novo Nordisk Foundation Center for Biosustainability, Technical University of Denmark, 2800 Kongens Lyngby, Denmark

\*Correspondence: e.marcellin@uq.edu.au

## ABSTRACT

Microbes that can recycle one-carbon (C1) greenhouse gases into fuels and chemicals are vital for the biosustainability of future industries. Acetogens are the most efficient known microbes for fixing carbon oxides CO<sub>2</sub> and CO. Understanding proteome allocation is important for metabolic engineering as it dictates metabolic fitness. Here, we use absolute proteomics to quantify intracellular concentrations for >1,000 proteins in the model-acetogen *Clostridium autoethanogenum* grown on three gas mixtures. We detect prioritisation of proteome allocation for C1 fixation and significant expression of proteins involved in the production of acetate and ethanol as well as proteins with unclear functions. The data also revealed which isoenzymes are important. Integration of proteomic and metabolic flux data demonstrated that enzymes catalyse high fluxes with high concentrations and high in vivo catalytic rates. We show that flux adjustments were dominantly accompanied with changing enzyme catalytic rates rather than concentrations. Our work serves as a reference dataset and advances systems-level understanding and engineering of acetogens.

## 29 INTRODUCTION

30 Increasing concerns about irreversible climate change are accelerating the shift to renewable, carbon-  
 31 free energy production (e.g., solar, wind, fuel cells). However, many fuels and chemicals will stay  
 32 carbon-based, and thus, technologies for their production using sustainable and renewable feedstocks  
 33 are needed to transition towards a circular bioeconomy. Moreover, the rising amount of solid waste  
 34 produced by human activities (e.g., municipal solid waste, lignocellulosic waste) will further endanger  
 35 our ecosystems' already critical state. Both challenges can be tackled by using organisms capable of  
 36 recycling gaseous one-carbon (C1) waste feedstocks (e.g., industrial waste gases [CO<sub>2</sub>, CO, CH<sub>4</sub>],  
 37 syngas from gasified biomass or municipal solid waste [CO, H<sub>2</sub>, CO<sub>2</sub>]) into fuels and chemicals at  
 38 industrial scale<sup>1-3</sup>.

39 As we transition into a new bioeconomy, a key feature of global biosustainability will be the  
 40 capacity to convert carbon oxides into products at industrial scale. Acetogens are the ideal biocatalysts  
 41 for this as they use the most energy-efficient pathway, the Wood-Ljungdahl pathway (WLP)<sup>4,5</sup>, for  
 42 fixing CO<sub>2</sub> into the central metabolite acetyl-CoA<sup>6-9</sup> and accept gas (CO, H<sub>2</sub>, CO<sub>2</sub>) as their sole carbon  
 43 and energy source<sup>5</sup>. Indeed, the model-acetogen *Clostridium autoethanogenum* is already being used  
 44 as a cell factory in industrial-scale gas fermentation<sup>3,10</sup>. The WLP is considered the first biochemical  
 45 pathway on Earth<sup>7,11-13</sup> and continues to play a critical role in the biogeochemical carbon cycle by  
 46 fixing an estimated 20% of the global CO<sub>2</sub><sup>6,14</sup>. While biochemical details of the WLP are well  
 47 described<sup>4,6,15</sup>, a quantitative understanding of acetogen metabolism is just emerging<sup>16,17</sup>. Notably,  
 48 recent systems-level analyses of acetogen metabolism have revealed mechanisms behind metabolic  
 49 shifts<sup>18-21</sup>, transcriptional architectures<sup>22,23</sup>, and features of translational regulation<sup>24,25</sup>. However, we  
 50 still lack an understanding of acetogen proteome allocation through the quantification of proteome-  
 51 wide intracellular protein concentrations. This fundamental knowledge is required for advancing  
 52 rational metabolic engineering of acetogen cell factories and for accurate in silico reconstruction of  
 53 their phenotypes using metabolic models<sup>1,2</sup>.

54 Quantitative description of an organism's proteome allocation through absolute proteome  
 55 quantification is valuable in several ways. Firstly, it enables us to understand prioritisation of the  
 56 energetically costly proteome resources among functional protein categories, metabolic pathways, and

single proteins<sup>26,27</sup>. This may also identify relevant proteins with unclear functions and high abundances. Secondly, some metabolic fluxes can be catalysed by isoenzymes and a comparison of their intracellular concentrations can indicate which are likely relevant in vivo and are thus targets for genetic perturbation experiments to validate in vivo functionalities<sup>28</sup>. Thirdly, integration of absolute proteomics and metabolic flux data enable the estimation of apparent in vivo catalytic rates of enzymes ( $k_{app}$ )<sup>26,29</sup>, which can be used to identify less-efficient enzymes as targets for improving pathways through metabolic and protein engineering. Absolute proteomics data also contribute to the curation of accurate genome-scale metabolic models.

Absolute proteome quantification is generally performed using label-free mass-spectrometry (MS) approaches without spike-in standards<sup>30,31</sup>. The major limitation of this approach is that accuracy of label-free estimated protein concentrations cannot be determined. Furthermore, the optimal model to convert MS signals (e.g., spectral counts, peak intensities) into protein concentrations remains unknown<sup>32–34</sup>. Label-based approaches using stable-isotope labelled (SIL) spike-ins of endogenous proteins are thus preferred for reliable absolute proteome quantification. This strategy relies on accurate absolute quantification of a limited set of intracellular proteins (i.e., anchors) using SIL spike-ins to establish a linear correlation between protein concentrations and their measured MS intensities<sup>32</sup>. Studies with the latter approach have determined a 1.5–2.4-fold error for label-free estimation of proteome-wide protein concentrations in multiple organisms<sup>28,35–41</sup>.

The aim of our work was to perform reliable absolute proteome quantification for the first time in an acetogen. We employed a label-based MS approach using SIL-protein spike-in standards to quantify SIL-based concentrations for 16 key proteins and label-free-based concentrations for >1,000 *C. autoethanogenum* proteins during autotrophic growth on three gas mixtures. This allowed us to explore global proteome allocation, uncover isoenzyme usage in central metabolism, and quantify regulatory principles associated with estimated  $k_{app}$ s. Our work provides an important reference dataset and advances the systems-level understanding and engineering of the ancient metabolism of acetogens.

83

## 84 RESULTS

85 **Absolute proteome quantification framework in the model-acetogen *C. autoethanogenum*.** We  
 86 performed absolute proteome quantification from autotrophic steady-state chemostat cultures of *C.*  
 87 *autoethanogenum* grown on three different gas mixtures: CO, syngas (CO+CO<sub>2</sub>+H<sub>2</sub>), or CO+H<sub>2</sub>  
 88 (termed “high-H<sub>2</sub> CO”) described before<sup>18,19</sup>. Briefly, four biological cultures of each gas mixture  
 89 were grown anaerobically on a chemically defined medium at 37 °C, pH of 5, and dilution rate ~1  
 90 day<sup>-1</sup> (specific growth rate ~0.04 h<sup>-1</sup>) without the use of heavy SIL substrates. The absolute proteome  
 91 quantification framework (Fig. 1) was built on using 19 synthetic heavy SIL-variants of key *C.*  
 92 *autoethanogenum* proteins covering central metabolism (Supplementary Table 1). The SIL-protein  
 93 standards were spiked in for quantification of intracellular concentrations of their endogenous light  
 94 counterparts. This framework ensures accurate absolute quantification compared to commonly used  
 95 peptide spike-ins. Spiking cell lysates with protein standards before sample clean-up and protein  
 96 digestion accounts for errors accompanying these critical steps<sup>30,31,42,43</sup>. Furthermore, selection of  
 97 peptides ensuring accurate absolute protein quantification without prior MS data is challenging as its  
 98 difficult to predict which peptides “fly” well<sup>30,31,42,43</sup>. In contrast, all proteotypic peptides from a  
 99 protein spike-in can be used for quantification.

100 We synthesised heavy-labelled lysine and arginine SIL-proteins using a cell-free wheat germ  
 101 extract platform as described previously<sup>18,44,45</sup> and quantified standard stocks using parallel reaction  
 102 monitoring (PRM) MS. Next, proteins were extracted from culture samples using an optimised  
 103 protocol maximising extraction yield<sup>18</sup> followed by spike-in of the 19 heavy SIL-proteins into light  
 104 cell lysates. We then used a data-independent acquisition (DIA) MS approach<sup>46</sup> to quantitate 1,243  
 105 proteins of *C. autoethanogenum* across 12 samples (quadruplicate cultures of three gas mixtures)  
 106 using a comprehensive spectral library consisting of whole-cell lysates, lysate fractions, and spike-in  
 107 SIL-proteins. Finally, we quantified intracellular concentrations for 16 key *C. autoethanogenum*  
 108 proteins using light-to-heavy ratios between endogenous and spike-in DIA MS intensities and further  
 109 used these 16 as anchor proteins for label-free estimation of ~1,043 protein concentrations through  
 110 establishing a linear correlation between protein concentrations and their measured MS intensities<sup>32</sup>.

We express intracellular protein concentrations in nanomoles of protein per gram of dry cell weight (nmol/gDCW).

**Absolute quantification of 16 anchor protein concentrations.** To ensure high confidence absolute quantification of anchor protein concentrations from the DIA MS data, we employed stringent criteria on top of the automated mProphet peak picking algorithm<sup>47</sup> within the software Skyline<sup>48</sup>. We also performed a dilution series experiment for each SIL-protein to increase accuracy (see Methods for details). Briefly, we kept only peaks with Gaussian shapes and without interference and precursors with highest Skyline quality metrics. Importantly, only peptides whose signal were above the lower limit of quantification (LLOQ) and within the linear dynamic quantification range in the dilution series experiment were used for anchor protein quantification (Supplementary Table 2). We thus used 106 high-confidence peptides for the absolute quantification of 16 anchor protein concentrations (Table 1; see also Fig. 5). High confidence of the intracellular concentrations for these key *C. autoethanogenum* proteins of central metabolism is supported both by the low average 11% coefficient of variation (CV) between biological quadruplicate cultures (Table 1) and the average 22% CV between different peptides of single proteins (Supplementary Table S2).

**Label-free estimation of proteome-wide protein concentrations.** Both high quality proteomics data and suitable anchor proteins are required for reliable label-free absolute proteome quantification. Our proteome-wide DIA MS data were highly reproducible with an average Pearson correlation coefficient of  $R = 0.99$  between biological replicates (Fig. 2a and Supplementary Fig. 1). We also found our anchor proteins suitable as their concentrations spanned across three orders of magnitude, and the summed mass accounted for  $\sim 1/3$  of the peptide mass injected into the mass spectrometer (Table 1 and Fig. 2b). We used the 16 anchor proteins (with 106 peptides) to determine the optimal label-free quantification model with the best linear fit between anchor protein concentrations and their measured DIA MS intensities using the aLFQ R package<sup>49</sup> as described before for SWATH MS<sup>28</sup> (Fig. 2c). Notably, we detected an average 1.5-fold cross-validated mean fold-error (CV-MFE; bootstrapping) for the label-free estimated anchor protein concentrations across samples (Fig. 2d).

The errors were distributed normally (Supplementary Fig. 2) with an average 95% CI of 0.3 (Fig. 2d). We then applied the optimal label-free quantification model to estimate ~1,043 protein concentrations in *C. autoethanogenum* (Supplementary Table 3).

Prior to the detailed analysis of proteome-wide protein concentrations, we further evaluated our label-free data accuracy beyond the 1.5-fold CV-MFE determined above. Firstly, the total proteome mass ( $1.2 \pm 0.1$   $\mu$ g; average  $\pm$  standard deviation) closely matched the 1  $\mu$ g peptide mass injected into the mass spectrometer (Fig. 2d). The data were also supported by a strong correlation between estimated protein concentrations and expected stoichiometries for equimolar (Fig. 3) and non-equimolar protein complexes (Supplementary Fig. 3). Notably, absolute protein concentrations of syngas cultures correlated well ( $R = 0.65$ ) with their respective absolute transcript expression levels determined before<sup>19</sup> (Supplementary Fig. 4). This result is similar to the correlations of absolute data seen in other steady-state cultures<sup>26,50</sup>. Altogether, we present the first absolute quantitative proteome dataset for a gas-fermenting acetogen that includes SIL-based concentrations for 16 key proteins and label-free estimates for over 1,000 *C. autoethanogenum* proteins during growth on three gas mixtures.

**C1 fixation dominates global proteome allocation.** Global proteome allocation amongst functional gene classifications was explored using proteomaps<sup>27</sup> and KEGG Orthology identifiers (KO IDs)<sup>51</sup>. The “treemap” structure defining the four-level hierarchy of our proteomaps (Supplementary Table 4) also included manually curated categories to accurately reflect acetogen metabolism (e.g., C1 fixation/WLP, Hydrogenases). As expected for autotrophic growth of an acetogen, the C1 fixation (Fig. 4) or WLP (Supplementary Fig. 5) categories dominated the proteome allocation with a  $\sim 1/3$  fraction, compared to Carbohydrate metabolism or Glycolysis/Gluconeogenesis. Notably, the data show that two genes—dihydrolipoamide dehydrogenase (LpdA; CAETHG\_RS07825) and glycine cleavage system H protein (GcvH; RS07795)—encoded by the WLP gene cluster were translated at very high levels (Fig. 4). This is important as both have unknown functions in *C. autoethanogenum* metabolism. Significant investment in expression of proteins involved in acetate and ethanol production (Supplementary Fig. 5) is consistent with  $1/3$ -to- $2/3$  of fixed carbon channelled into these two growth by-products across the three gas mixtures<sup>18,19</sup>. The 11% proteome fraction of category

Translation (Fig. 4) is expected for cells growing at a specific growth rate  $\sim 0.04 \text{ h}^{-1}$  based on absolute proteomics data from *Escherichia coli*<sup>26,39,52</sup>. The notable proteome allocation for Amino acid metabolism and particularly the high abundance of ketol-acid reductoisomerase (IlvC; RS00580) are surprising since metabolic fluxes through 2,3-butanediol and branched-chain amino acid pathways were low under these growth conditions<sup>18,19</sup>. In addition, numerous proteins with unknown or unclear functions (coloured grey in proteomaps) are highly expressed (e.g., RS12590, RS08610, RS08145), highlighting the need for global mapping of genotype-phenotype relationships in acetogens. In general, proteome allocation was highly similar between the three gas mixtures (Supplementary Table 3). This result is unsurprising given the few relative protein expression differences detected previously<sup>18</sup>.

**Enzyme usage revealed in central metabolism.** Next, we focused on uncovering enzyme usage in acetogen central metabolism (Fig. 5). This contains enzymes of the WLP, acetate, ethanol, and 2,3-butanediol production pathways, hydrogenases, and the Nfn transhydrogenase, which together carry >90% of the carbon and most of the redox flow in *C. autoethanogenum*<sup>18–20</sup>. Multiple metabolic fluxes in these pathways can be catalysed by isoenzymes and absolute proteomics data can indicate which of the isoenzymes are likely relevant in vivo. While the carbon monoxide dehydrogenase (CODH) AcsA (RS07861 □ 62) that forms the bifunctional CODH/ACS complex with the acetyl-CoA synthase<sup>53</sup> (AcsB; RS07800) is essential for *C. autoethanogenum* growth on gas as confirmed in mutagenesis studies<sup>54</sup>, the higher concentrations of the dispensable monofunctional CODH CooS1 (RS14775) suggest it may also play a role in CO oxidation (Fig. 4, 5), in addition to CO<sub>2</sub> reduction<sup>54</sup>. Additionally, our proteomics data show high abundance of the primary acetaldehyde:ferredoxin oxidoreductase (AOR1; RS00440) and this support the emerging understanding that in *C. autoethanogenum* ethanol is dominantly produced using the AOR1 activity via acetate, instead of directly from acetyl-CoA via acetaldehyde using mono- or bifunctional activities<sup>18,19,55,56</sup> (Fig. 5). Furthermore, the data suggest that the specific alcohol dehydrogenase (Adh4; RS08920) is responsible for reducing acetaldehyde to ethanol, a key reaction in terms of carbon and redox metabolism. The high abundance of the electron-bifurcating hydrogenase HytA-E complex (RS13745 □ 70) compared to alternative hydrogenases confirms that it is the main H<sub>2</sub>-oxidiser<sup>57,58</sup> (Fig. 5). This is consistent with

the fact that in the presence of H<sub>2</sub> all the CO<sub>2</sub> fixed by the WLP is reduced to formate using H<sub>2</sub> by the HytA-E and formate dehydrogenase (FdhA; RS13725) enzyme complex activity<sup>18,19</sup>. Despite the proteomics evidence, genetic perturbations are required to determine condition-specific in vivo functionalities of isoenzymes in acetogens unequivocally.

The overall most abundant protein was the formate-tetrahydrofolate ligase (Fhs; RS07850), a key enzyme in the WLP (Fig. 5, 4). Despite the high abundance, its expression might still be rate-limiting (see below). Another key enzyme for acetogens is AcsB because of its essentiality for acetyl-CoA synthesis by the CODH/ACS complex. AcsB is linked to the WLP by the corrinoid iron sulfur proteins AcsC (RS07810) and AcsD (RS07815) that supply the methyl group to AcsB. Interestingly, the ratio of AcsCD-to-AcsB increased from 1.7 (CO) to 2.3 (syngas) to 2.9 (high-H<sub>2</sub> CO), suggesting that the primary role of the CODH/ACS complex shifted from CO oxidation towards acetyl-CoA synthesis, likely because increased H<sub>2</sub> uptake could replace the supply of reduced ferredoxin from CO oxidation. Concurrently, the Nfn transhydrogenase (RS07665) levels that act as a redox valve in acetogens<sup>20</sup> are maintained high (Fig. 5), potentially to rapidly respond to redox perturbations. We conclude that absolute quantitative proteomics can significantly contribute to a systems-level understanding of metabolism, particularly in less-studied organisms.

**Integration of absolute proteomics and flux data yields in vivo enzyme catalytic rates.** Absolute proteomics data enable estimation of intracellular catalytic working rates of enzymes when metabolic flux rates are known<sup>26,29</sup>. We thus calculated apparent in vivo catalytic rates of enzymes, denoted as  $k_{app}$  (s<sup>-1</sup>)<sup>26</sup>, as the ratio of specific flux rate (mmol/gDCW/h) determined before<sup>18</sup> and protein concentration (nmol/gDCW) (see Methods). This produced  $k_{app}$  values for 13 and 48 enzymes/complexes using either anchor or label-free protein concentrations, respectively (Supplementary Table 5 and Fig. 5, 6). The first two critical steps for carbon fixation in the methyl branch of the WLP (i.e., CO to formate) are catalysed at high rates (Fig. 5). Notably, FdhA showed a  $k_{app}$  ~30 s<sup>-1</sup> for CO<sub>2</sub> reduction without H<sub>2</sub> during growth on CO only, which is similar to in vitro  $k_{cat}$  data of formate dehydrogenases in other acetogens<sup>59,60</sup>. Interestingly, the next step of formate reduction was catalysed potentially by a less-efficient enzyme—Fhs—as its  $k_{app}$  of ~3 s<sup>-1</sup> is significantly



lower compared to other WLP enzymes (Fig. 5). Overall, enzymes catalysing reactions in high flux pathways such as the WLP and acetate and ethanol production have higher  $k_{app}$ s than those of downstream from conversion of acetyl-CoA to pyruvate (Fig. 5). Indeed, enzymes catalysing high metabolic fluxes in *C. autoethanogenum* have both higher concentrations and higher catalytic rates compared to enzymes catalysing lower fluxes as both specific flux rates and enzyme concentrations (Kendall's  $\tau = 0.56$ , p-value =  $5 \times 10^{-9}$ ) and flux and  $k_{app}$  ( $\tau = 0.45$ , p-value =  $2 \times 10^{-6}$ ) were significantly correlated (Fig. 6a), as seen before for other organisms<sup>26,61</sup>.

Having acquired absolute proteomics data for *C. autoethanogenum* growth on three gas mixtures with different metabolic flux profiles also allowed us to determine the impact of change in enzyme concentration and its catalytic rate for adjusting metabolic flux rates. Two extreme examples are the reactions catalysed by the HytA-E (Fig. 6b) and the Nfn (Fig. 5) complexes where flux adjustments were accompanied with large changes in  $k_{app}$ s rather than in enzyme concentrations. Flux changes in high flux pathways such as the WLP and acetate and ethanol production also coincided mainly with  $k_{app}$  changes (Fig. 6b). This principle seems to be dominant in *C. autoethanogenum* as 90% of flux changes were not regulated through enzyme concentrations (i.e., post-translational regulation; Supplementary Table 6) when comparing all statistically significant flux changes between the three gas mixtures with respective enzyme expression changes (see Methods).

## DISCUSSION

The looming danger of irreversible climate change and harmful effects of solid waste accumulation are pushing humanity to develop and adopt sustainable technologies for renewable production of fuels and chemicals and for waste recycling. Acetogen gas fermentation offers great potential to tackle both challenges through recycling waste feedstocks (e.g., industrial waste gases, gasified biomass or municipal solid waste) into fuels and chemicals<sup>1,2</sup>. Although the quantitative understanding of acetogen metabolism has recently improved<sup>16,17</sup>, a quantitative description of acetogen proteome allocation was missing. This is needed to advance their metabolic engineering into superior cell factories and accurate in silico reconstruction of their phenotypes<sup>1,2</sup>. Thus, we performed absolute

proteome quantification in the model-acetogen *C. autoethanogenum* grown autotrophically on three gas mixtures.

Our absolute proteome quantification framework relied on SIL-protein spike-in standards and DIA MS analysis to ensure high confidence of the determined intracellular concentrations for 16 key *C. autoethanogenum* proteins. We further used these proteins as anchor proteins for label-free estimation of >1,000 protein concentrations. This enabled us to determine the optimal label-free quantification model for our data to infer protein concentrations from MS intensities, which remains unknown in common label-free approaches not utilising spike-in standards<sup>32,33</sup>. More importantly, label-free estimated protein concentrations using the latter approach are questionable as their accuracy cannot be determined. We determined an excellent average error of 1.5-fold for our label-free estimated protein concentrations based on 16 anchor proteins and a bootstrapping approach. This error is in the same range as described in previous studies using SIL spike-in standards for absolute proteome quantification<sup>28,35–41</sup>. Further, we also observed a good match both between estimated and injected proteome mass into the mass spectrometer and between protein concentrations and expected protein complex stoichiometries. We conclude that label-free estimation of proteome-wide protein concentrations using SIL-protein spike-ins and state-of-the-art MS analysis is reasonably accurate.

Quantification of acetogen proteome allocation during autotrophic growth expectedly showed prioritisation of proteome resources for fixing carbon through the WLP, in line with transcript expression data in *C. autoethanogenum*<sup>19,56</sup>. The allocation of one third of the total proteome for C1 fixation is higher than proteome allocation for carbon fixation through glycolysis during heterotrophic growth of other microorganisms<sup>26,36</sup>. High abundances of other key enzymes of acetogen central metabolism were also expected as the WLP, acetate and ethanol production pathways, hydrogenases, and the Nfn transhydrogenase carry >90% of the carbon and most of the redox flow in *C. autoethanogenum*<sup>18–20</sup>. However, very high expression of the two genes—LpdA and GcvH—of the WLP gene cluster with unknown functions in *C. autoethanogenum* is striking, raising the question whether their function in *C. autoethanogenum* could also be to link WLP and glycine synthase-reductase pathways, as recently proposed for another acetogen<sup>62</sup>. Since many other proteins with unknown or

unclear functions were also highly abundant, global mapping of genotype-phenotype relationships in acetogens is much needed.

The in vivo functionalities of isoenzymes are not clear for multiple key metabolic fluxes in acetogen central metabolism and absolute proteomics data can indicate which isoenzymes are likely relevant. Oxidation of CO or reduction of CO<sub>2</sub> is a fundamental step for all acetogens and known to be catalysed by three CODHs in *C. autoethanogenum*<sup>54</sup>. Though only AcsA that forms the bifunctional CODH/ACS complex with the acetyl-CoA synthase<sup>53</sup> is essential for growth on gas<sup>54</sup>, we detected higher concentrations of the monofunctional CODH CooS1, which deletion strain shows intriguing phenotypes<sup>54</sup>. Concurrently, our data suggest that prioritisation of CODH/ACS activity between CO oxidation and acetyl-CoA synthesis is sensitive to H<sub>2</sub> availability. Thus, further studies are required to decipher condition-dependent functionalities of CODHs. In addition to CODHs, the biochemical understanding of ethanol production is important in terms of both carbon and redox metabolism. Our data confirm that in *C. autoethanogenum* ethanol is predominantly produced via acetate by AOR1<sup>18,19,55,56</sup> and more importantly, indicate for the first time that AOR1 activity is followed by Adh4 (previously characterised as butanol dehydrogenase<sup>63</sup>) for reduction of acetaldehyde to ethanol. These observations call for large-scale genetic perturbation experiments to determine unequivocally the condition-specific in vivo functionalities of isoenzymes in acetogens.

Absolute proteomics data offer a unique opportunity to estimate apparent in vivo catalytic rates of enzymes ( $k_{app}$ )<sup>26,29</sup> if also metabolic flux data are available. These data are particularly valuable for more accurate in silico reconstruction of phenotypes using protein-constrained genome-scale metabolic models<sup>64,65</sup>. While in vitro  $k_{cat}$  and in vivo  $k_{app}$  data generally correlate<sup>29</sup>, models using maximal  $k_{app}$  values show better prediction of protein abundances<sup>66</sup>. Furthermore, information of  $k_{app}$ s can infer less-efficient enzymes as targets for improving pathways through metabolic and protein engineering. For example, protein engineering of Fhs (catalysing formate reduction) might improve WLP throughput and carbon fixation since its  $k_{app}$  was significantly lower compared to other pathway enzymes. At the same time, the large change in  $k_{app}$ s of the abundant electron-bifurcating hydrogenase HytA-E and the Nfn transhydrogenase complexes indicate capacity for the cells to rapidly respond to H<sub>2</sub> availability and redox perturbations, which may be critical for metabolic robustness of acetogens<sup>20</sup>.

Overall, we detected both higher concentrations and  $k_{app}$ s for enzymes catalysing higher metabolic fluxes, which is believed to arise from an evolutionary push towards reducing protein production costs for enzymes carrying high flux<sup>61</sup>. The observation that 90% of flux changes in *C. autoethanogenum* were not regulated through changes in enzyme concentrations is not surprising for a metabolism that operates at the thermodynamic edge of feasibility<sup>16,17</sup> since post-translational regulation of fluxes is energetically least costly. Further research is needed to identify which mechanism from post-translational protein modification, allosteric regulation, or substrate concentration change is responsible for post-translational regulation of fluxes.

We have produced the first absolute proteome quantification in an acetogen and thus provided understanding of global proteome allocation, isoenzyme usage in central metabolism, and regulatory principles of in vivo enzyme catalytic rates. This fundamental knowledge has potential to advance both rational metabolic engineering of acetogen cell factories and accurate in silico reconstruction of their phenotypes<sup>1,2,64,65</sup>. Our study also highlights the need for large-scale mapping of genotype-phenotype relationships in acetogens to infer in vivo functionalities of isoenzymes and proteins with unknown or unclear functions. This absolute proteomics dataset serves as a reference towards a better systems-level understanding of the ancient metabolism of acetogens.

## METHODS

**Bacterial strain and culture growth conditions.** Absolute proteome quantification was performed from high biomass concentration (~1.4 gDCW/L) steady-state autotrophic chemostat cultures of *C. autoethanogenum* growing on three different gas mixtures with culturing conditions described in our previous works<sup>18,19</sup>. Briefly, four biological replicate chemostat cultures of *C. autoethanogenum* strain DSM 19630 were grown on a chemically defined medium (without yeast extract) either on CO (~60% CO and 40% Ar), syngas (~50% CO, 20% H<sub>2</sub>, 20% CO<sub>2</sub>, and 10% N<sub>2</sub>/Ar), or CO+H<sub>2</sub>, termed “high-H<sub>2</sub> CO”, (~15% CO, 45% H<sub>2</sub>, and 40% Ar) under strictly anaerobic conditions. The bioreactors were maintained at 37 °C, pH of 5, and dilution rate ~1 day<sup>-1</sup> (specific growth rate ~0.04 h<sup>-1</sup>).

**Cell-free synthesis of stable-isotope labelled protein standards.** Twenty proteins covering *C. autoethanogenum* central carbon metabolism, the HytA-E hydrogenase, and a ribosomal protein (Supplementary Table 1) were selected for cell-free synthesis of SIL-proteins as described in ref.<sup>18</sup>. Briefly, genes encoding for these proteins were synthesised by commercial gene synthesis services (Biomatik). Target genes were sub-cloned into the cell-free expression vector pEUE01-His-N2 (CellFree Sciences) and transformed into *Escherichia coli* DH5α from which plasmid DNA was extracted and purified. Correct gene insertion into the pEUE01-His-N2 was verified by DNA sequencing. Subsequently, cell-free synthesis of His-tag fused *C. autoethanogenum* proteins was performed using the bilayer reaction method with the wheat germ extract WEPRO8240H (CellFree Sciences) as described previously<sup>44,45</sup>. mRNAs for cell-free synthesis were prepared by an in vitro transcription reaction while in vitro translation of target proteins was performed using a bilayer reaction where the translation layer was supplemented with L-Arg-<sup>13</sup>C<sub>6</sub>, <sup>15</sup>N<sub>4</sub> and L-Lys-<sup>13</sup>C<sub>6</sub>, <sup>15</sup>N<sub>2</sub> (Wako) at final concentrations of 20 mM to achieve high efficiency (>99 %) for stable-isotope labelling of proteins. The in vitro synthesised SIL-protein sequences also contained an N-terminal amino acid sequence GYSFTTTAEK that was later used as a tag for quantification of the SIL-protein stock concentration. Subsequently, SIL-proteins were purified using the Ni-Sepharose High-Performance resin (GE Healthcare Life Sciences) and precipitated using methanol:chloroform:water precipitation in Eppendorf Protein LoBind<sup>®</sup> tubes. Lastly, precipitated SIL-proteins were reconstituted in 104 µL of 8 M urea ([UA]; Sigma-Aldrich) in 0.1 M Trizma<sup>®</sup> base (pH 8.5) by vigorous vortexing and stored at -80 °C until further use.

**Absolute quantification of SIL-protein standards using PRM MS.** Concentrations of the twenty synthesised SIL-protein standard stocks were determined using PRM MS preceded by in-solution digestion of proteins and sample desalting and preparation for MS analysis.

### ***Sample preparation***

Only Eppendorf Protein LoBind<sup>®</sup> tubes and pipette tips were used for all sample preparation steps. Firstly, 20 µL of UA was added to 4 µL of the SIL-protein standard stock used to determine the stock concentration, and the mix was vortexed. Then, 1 µL of 0.2 M DTT (Promega) was added, followed

by vortexing and incubation for 1 h at 37 °C to reduce disulphide bonds. Sulfhydryl groups were alkylated with 2 µL of 0.5 M iodoacetamide (IAA; Sigma-Aldrich), vigorous vortexing, and incubation for 30 min at room temperature in the dark. Next, 75 µL of 25 mM ammonium bicarbonate was added to dilute UA down to 2 M concentration. Subsequently, 2 pmol (2 µL of stock) of the non-labelled AQUA<sup>®</sup> peptide HLEAAKGYSFTTTAEKAAELHK (Sigma-Aldrich) containing the quantification tag sequence GYSFTTTAEK was added to enable quantification of SIL-protein stock concentrations using MS analysis based on the ratio of heavy-to-light GYSFTTTAEK signals (see below). Protein digestion was performed for 16 h at 37 °C with 0.1 µg of Trypsin/Lys-C mix (1 µL of stock; Promega) and stopped by lowering pH to 3 by the addition of 5 µL of 10% (v/v) trifluoroacetic acid (TFA).

Samples were desalted using C<sub>18</sub> ZipTips (Merck Millipore) as follows: the column was wetted using 0.1% (v/v) formic acid (FA) in 100% acetonitrile (ACN), equilibrated with 0.1% FA in 70% (v/v) ACN, and washed with 0.1% FA before loading the sample and washing again with 0.1% FA. Peptides were eluted from the ZipTips with 0.1% FA in 70% ACN. Finally, samples were dried using a vacuum-centrifuge (Eppendorf) at 30 °C until dryness followed by reconstitution in 12 µL of 0.1% FA in 5% ACN for subsequent MS analysis.

#### ***LC method for PRM MS***

A Thermo Fisher Scientific UltiMate 3000 RSLCnano UHPLC system was used to elute the samples. Each sample was initially injected (6 µL) onto a Thermo Fisher Acclaim PepMap C<sub>18</sub> trap reversed-phase column (300 µm x 5 mm nano viper, 5 µm particle size) at a flow rate of 15 µL/min using 2% ACN for 3 min with the solvent going to waste. The trap column was switched in-line with the separation column (GRACE Vydac Everest C<sub>18</sub>, 300Å 150 µm x 150 mm, 2 µm) and the peptides were eluted using a flowrate of 3 µL/min using 0.1% FA in water (buffer A) and 80% ACN in buffer A (buffer B) as mobile phases for gradient elution. Following 3 min isocratic of 3% buffer B, peptide elution employed a 3-40% ACN gradient for 28 min followed by 40-95% ACN for 1.5 min and 95% ACN for 1.5 min at 40 °C. The total elution time was 50 min including a 95% ACN wash and a re-equilibration step.

#### ***PRM MS data acquisition***

The eluted peptides from the C18 column were introduced to the MS via a nano-ESI and analysed using the Thermo Fisher Scientific Q-Exactive HF-X mass spectrometer. The electrospray voltage was 1.8 kV in positive ion mode, and the ion transfer tube temperature was 250 °C. Full MS-scans were acquired in the Orbitrap mass analyser over the range  $m/z$  550–560 with a mass resolution of 30,000 (at  $m/z$  200). The AGC target value was set at 1.00E+06 and maximum accumulation time 50 ms for full MS-scans. The PRM inclusion list included two mass values of 552.7640 and 556.7711. MS/MS spectra were acquired in the Orbitrap mass analyser with a mass resolution of 15,000 (at  $m/z$  200). The AGC target value was set at 1.00E+06 and maximum accumulation time 30 ms for MS/MS with an isolation window of 2  $m/z$ . The loop count was set at 14 to gain greater MS/MS data. Raw PRM MS data have been deposited to Panorama at [https://panoramaweb.org/Valgepea\\_Cauto\\_PRM.url](https://panoramaweb.org/Valgepea_Cauto_PRM.url) (private reviewer account details: username: panorama+reviewer27@proteinms.net; password: hSAwNUAL) with a ProteomeXchange Consortium (<http://proteomecentral.proteomexchange.org>) dataset identifier PXD025760.

### ***PRM MS data analysis***

Analysis of PRM MS data was performed using the software Skyline<sup>48</sup>. The following parameters were used to extract PRM MS data for the quantification tag sequence GYSFTTTAEK: three precursor isotope peaks with a charge of 2 (++) were included (monoisotopic; M+1; M+2); five of the most intense y product ions from ion 3 to last ion of charge state 1 and 2 among the precursor were picked; chromatograms were extracted with an ion match mass tolerance of 0.05  $m/z$  for product ions by including all matching scans; full trypsin specificity with two missed cleavages allowed for peptides with a length of 8-25 AAs; cysteine carbamidomethylation as a fixed peptide modification. Additionally, peptide modifications included heavy labels for lysine and arginine as <sup>13</sup>C(6)<sup>15</sup>N(2)/+8.014 Da (K) and <sup>13</sup>C(6)<sup>15</sup>N(4)/+10.008 Da (R), respectively. This translated into the SIL-proteins and the non-labelled AQUA<sup>®</sup> peptide possessing the tag GYSFTTTAEK with  $m/z$  of 556.7711 and 552.7640, respectively. Hence, the concentrations of SIL-protein stocks were calculated based on the ratio of heavy-to-light GYSFTTTAEK signals and the spike-in of 2 pmol of the non-labelled AQUA<sup>®</sup> peptide (see above). High accuracy of quantification was evidenced by the very high similarity between both precursor peak areas and expected isotope distribution ( $R^2 > 0.99$ ; idotp in

Skyline) and heavy and light peak areas ( $R^2 > 0.99$ ; rdotp in Skyline) for all SIL-protein standard stocks. No heavy GYSFTTTAEK signal was detected for protein CAETHG\_RS16140 (an acetylating acetaldehyde dehydrogenase in the NCBI annotation of sequence NC\_022592.1<sup>67</sup>), thus 19 SIL-proteins could be used for following absolute proteome quantification in *C. autoethanogenum* (Supplementary Table 1). PRM MS data with all Skyline processing settings can be viewed and downloaded from Panorama at [https://panoramaweb.org/Valgepea\\_Cauto\\_PRM.url](https://panoramaweb.org/Valgepea_Cauto_PRM.url) (private reviewer account details: username: panorama+reviewer27@proteinms.net; password: hSAwNUAL).

**Absolute proteome quantification in *C. autoethanogenum* using DIA MS.** We used 19 synthetic heavy SIL variants (see above) of key *C. autoethanogenum* proteins (Supplementary Table 1) as spike-in standards for quantification of intracellular concentrations of their non-SIL counterparts using a DIA MS approach<sup>46</sup>. Also, we performed a dilution series experiment for the spike-in SIL-proteins to ensure accurate absolute quantification. We refer to these 19 intracellular proteins as anchor proteins that were further used to estimate proteome-wide absolute protein concentrations in *C. autoethanogenum*. This was achieved by determining the best linear fit between anchor protein concentrations and their measured DIA MS intensities using the same strategy as described previously<sup>28</sup>.

#### ***Preparation of spike-in SIL-protein standard mix and dilution series samples***

Only Eppendorf Protein LoBind<sup>®</sup> tubes and pipette tips were used for all preparation steps. The 19 spike-in SIL-protein standards that could be used for absolute proteome quantification in *C. autoethanogenum* (see above) were mixed in two lots: 1) ‘sample spike-in standard mix’: SIL-protein quantities matching estimated intracellular anchor protein quantities (i.e., expected light-to-heavy [L/H] ratios of ~1) based on label-free absolute quantification of the same samples in our previous work<sup>18</sup>; 2) ‘dilution series standard mix’: SIL-protein quantities doubling the estimated intracellular anchor protein quantities for the dilution series sample with the highest SIL-protein concentrations.

To ensure accurate absolute quantification of anchor protein concentrations, a dilution series experiment was performed to determine the linear dynamic quantification range and LLOQ for each of the 19 spike-in SIL-proteins. Dilution series samples were prepared by making nine 2-fold dilutions



of the ‘dilution series spike-in standard mix’ (i.e., 10 samples total for dilution series with a 512-fold concentration span) in a constant *C. autoethanogenum* cell lysate background (0.07 µg/µL; 10 µg/tube) serving as a blocking agent to avoid loss of purified SIL-proteins (to container and pipette tip walls) and as a background proteome for accurate MS quantification of the linear range and LLOQ for anchor proteins.

### **Sample preparation**

*C. autoethanogenum* cultures were sampled for proteomics by immediate pelleting of 2 mL of culture using centrifugation (25,000 × *g* for 1 min at 4 °C) and stored at -80 °C until analysis. Details of protein extraction and protein quantification in cell lysates are described previously<sup>18</sup>. In short, thawed cell pellets were suspended in lysis buffer (containing SDS, DTT, and Trizma® base) and cell lysis was performed using glass beads and repeating a ‘lysis cycle’ consisting of heating, bead beating, centrifugation, and vortexing before protein quantification using the 2D Quant Kit (GE Healthcare Life Sciences).

Sample preparation and protein digestion for MS analysis was based on the filter-aided sample preparation (FASP) protocol<sup>68</sup>. The following starting material was loaded onto an Amicon® Ultra-0.5 mL centrifugal filter unit (nominal molecular weight cut-off of 30,000; Merck Millipore): 1) 50 µg of protein for one culture sample from each gas mixture (CO, syngas, or high-H<sub>2</sub> CO) for building the spectral library for DIA MS data analysis (samples 1-3); 2) 7 µg of protein for one culture sample from either syngas or high-H<sub>2</sub> CO plus ‘sample spike-in standard mix’ for including spike-in SIL-protein data to the spectral library (samples 4-5); 3) 15 µg of protein for all 12 culture samples (biological quadruplicates from CO, syngas, and high-H<sub>2</sub> CO) plus ‘sample spike-in standard mix’ for performing absolute proteome quantification in *C. autoethanogenum* (samples 6-17); 4) ten dilution series samples with 10-15 µg of total protein (*C. autoethanogenum* cell lysate background plus ‘dilution series spike-in standard mix’) for performing the dilution series experiment for the 19 spike-in SIL-proteins (see above) (samples 18-27).

Samples containing SIL-proteins (samples 4-27) were incubated at 37 °C for 1 h to reduce SIL-protein disulphide bonds (cell lysate contained DTT). Details of the FASP workflow are described before<sup>18</sup>. In short, samples were washed with UA, sulfhydryl groups alkylated with IAA,

proteins digested using a Trypsin/Lys-C mix, and peptides eluted from the filter with 60  $\mu$ L of ammonium bicarbonate. Next, 50  $\mu$ L of samples 1-3 were withdrawn and pooled for performing high pH reverse-phase fractionation as described previously<sup>18</sup> for expanding the spectral library for DIA MS data analysis, yielding eight fractions (samples 28-35). Subsequently, all samples were vacuum-centrifuged at 30 °C until dryness followed by reconstitution of samples 1-3 and 4-35 in 51 and 13  $\mu$ L of 0.1% FA in 5% ACN, respectively. Finally, total peptide concentration in each sample was determined using the Pierce<sup>TM</sup> Quantitative Fluorometric Peptide Assay (Thermo Fisher Scientific) to ensure that the same total peptide amount across samples 1-17 and 28-35 (excluding samples 18-27, see below) could be injected for DIA MS analysis.

### ***LC method for data-dependent acquisition (DDA) and DIA MS***

Details of the LC method employed for generating the spectral library using DDA and for DIA sample runs are described previously<sup>20</sup>. In short, a Thermo Fisher Scientific UHPLC system including C<sub>18</sub> trap and separation columns was used to elute peptides with a gradient and total elution time of 110 min. For each DDA and DIA sample run, 1  $\mu$ g of peptide material from protein digestion was injected, except for dilution series samples (samples 18-27 above) that were injected in a constant volume of 3  $\mu$ L to maintain the dilution levels of the ‘dilution series spike-in standard mix’.

### ***DDA MS spectral library generation***

The following 13 samples were analysed on the Q-Exactive HF-X in DDA mode to yield the spectral library for DIA MS data analysis: 1) three replicates of one culture sample from each gas mixture (CO, syngas, or high-H<sub>2</sub> CO) (samples 1-3 above); 2) three replicates of one culture sample from either syngas or high-H<sub>2</sub> CO plus ‘sample spike-in standard mix’ (samples 4-5); 3) eight high pH reverse-phase fractions of a pool of samples from each gas mixture (samples 28-35).

Details of DDA MS acquisition for generating the spectral library are described before<sup>20</sup>. In short, eluted peptides from the C<sub>18</sub> column were introduced to the MS *via* a nano-ESI and analysed using the Q-Exactive HF-X with an Orbitrap mass analyser. The DDA MS spectral library for DIA MS data confirmation and quantification using the software Skyline<sup>48</sup> was created using the Proteome Discoverer 2.2 software (Thermo Fisher Scientific) and its SEQUEST HT search as described previously<sup>18</sup>. The final .pd result file contained peptide-spectrum matches (PSMs) with q-values

estimated at 1% false discovery rate (FDR) for peptides  $\geq 4$  AAs. The generated spectral library file has been deposited to the ProteomeXchange Consortium (<http://proteomecentral.proteomexchange.org>) via the PRIDE partner repository<sup>69</sup> with the dataset identifier PXD025732 (private reviewer account details: username: reviewer\_pxd025732@ebi.ac.uk; password: QdkJHXy0).

### ***DIA MS data acquisition***

Details of DIA MS acquisition are described before<sup>20</sup>. In short, as for DDA MS acquisition, eluted peptides were introduced to the MS via a nano-ESI and analysed using the Q-Exactive HF-X with an Orbitrap mass analyser. DIA was achieved using an inclusion list: m/z 395–1100 in steps of 15 amu and scans cycled through the list of 48 isolation windows with a loop count of 48. In total, DIA MS data was acquired for 22 samples (samples 6-27 defined in section *Sample preparation*). Raw DIA MS data have been deposited to the ProteomeXchange Consortium (<http://proteomecentral.proteomexchange.org>) via the PRIDE partner repository<sup>69</sup> with the dataset identifier PXD025732 (private reviewer account details: username: reviewer\_pxd025732@ebi.ac.uk; password: QdkJHXy0).

### ***DIA MS data analysis***

DIA MS data analysis was performed with Skyline<sup>48</sup> as described before<sup>18</sup> with the following modifications: 1) 12 manually picked high confidence endogenous peptides present in all samples and spanning the elution gradient were used for iRT alignment through building an RT predictor; 2) outlier peptides from iRT regression were removed; 3) a minimum of three isotope peaks were required for a precursor; 4) single peptide per spike-in SIL-protein was allowed for anchor protein absolute quantification while at least two peptides per protein were required for label-free estimation of proteome-wide protein concentrations; 5) extracted ion chromatograms (XICs) were transformed using Savitzsky-Golay smoothing. Briefly, the .pd result file from Proteome Discoverer was used to build the DIA MS spectral library and the mProphet peak picking algorithm<sup>47</sup> within Skyline was used to separate true from false positive peak groups (per sample) and only peak groups with q-value<0.01 (representing 1% FDR) were used for further quantification. We confidently quantitated 7,288 peptides and 1,243 proteins across all samples and 4,887 peptides and 1,043 proteins on average

within each sample for estimating proteome-wide absolute protein concentrations. For absolute quantification of anchor protein concentrations, we additionally manually: 1) removed integration of peaks showing non-Gaussian shapes or interference from other peaks; 2) removed precursors with similarity measures of  $R^2 < 0.9$  between product peak areas and corresponding intensities in the spectral library (dotp in Skyline), precursor peak areas and expected isotope distribution (idotp), or heavy and light peak areas (rdotp). After analysis in Skyline, 17 spike-in SIL-proteins remained for further analysis as protein CAETHG\_RS14410 was not identified in DIA MS data while CAETHG\_RS18395 did not pass quantification filters (Supplementary Table 1). DIA MS data with all Skyline processing settings can be viewed and downloaded from Panorama at [https://panoramaweb.org/Valgepea\\_Cauto\\_Anchors.url](https://panoramaweb.org/Valgepea_Cauto_Anchors.url) for anchor protein absolute quantification and at [https://panoramaweb.org/Valgepea\\_Cauto\\_LF.url](https://panoramaweb.org/Valgepea_Cauto_LF.url) for estimating proteome-wide absolute protein concentrations (private reviewer account details: username: panorama+reviewer27@proteinms.net; password: hSAwNUAL).

#### ***Absolute quantification of anchor protein concentrations***

We employed further stringent criteria on top of the output from Skyline analysis to ensure high confidence absolute quantification of 17 anchor protein concentrations. Firstly, precursor with highest heavy intensity for the highest ‘dilution series spike-in standard mix’ sample in the dilution series (DS01) was kept while others were deleted. Peptides quantified in less than three biological replicates within a gas mixture, with no heavy signal for DS01 sample, or with heavy signals for less than three continuous dilution series samples were removed. Next, we utilised the dilution series experiment to only keep signals over the LLOQ and within the linear dynamic quantification range. For this, correlation between experimental and expected peptide L/H signal ratios for each peptide across the dilution series was made to determine the LLOQ and calculate correlation, slope, and intercept between MS signal and spike-in level (Supplementary Table 2). Only peptides showing correlation  $R^2 > 0.95$ ,  $0.95 < \text{slope} < 1.05$ , and  $-0.1 < \text{intercept} < 0.1$  for the dilution series were kept. This ensured that we were only using peptides within the linear dynamic range. The remaining peptides were further filtered for each culture sample by removing peptides whose light or heavy signal was below the LLOQ in the dilution series. Subsequently, only peptides were kept with L/H ratios for at

least three biological replicates cultures for each gas mixture (i.e.,  $\geq 9$  data points). Finally, we aimed to detect outlier peptides by calculating the % difference of a peptide's L/H ratio from the average L/H ratio of all peptides for a given protein for every sample. Peptides were considered outliers and thus removed if the average difference across all samples was  $>50\%$  or if the average difference within biological replicate cultures was  $>50\%$ . After the previous stringent criteria were applied, 106 high-confidence peptides remained (Supplementary Table 2) for the quantification of 16 anchor protein concentrations since CAETHG\_RS01830 was lost during manual analysis (Table 1 and Supplementary Table 1). Proteins CAETHG\_RS13725 and CAETHG\_RS07840 were excluded from the high- $H_2$  CO culture dataset as their calculated concentrations varied  $>50\%$  between biological replicates. Data of one high- $H_2$  CO culture was excluded from further analysis due to difference from bio-replicates likely due to challenges with MS analysis.

#### ***Label-free estimation of proteome-wide protein concentrations***

We used the anchor proteins to estimate proteome-wide protein concentrations in *C. autoethanogenum* by determining the best linear fit between anchor protein concentrations and their measured DIA MS intensities using the aLFQ R package<sup>49</sup> and the same strategy as described for SWATH MS<sup>28</sup>. Briefly, aLFQ used anchor proteins and cross-validated model selection by bootstrapping to determine the optimal model within various label-free absolute proteome quantification approaches (e.g., TopN, iBAQ). The approach can obtain the model with the smallest error between anchor protein concentrations determined using SIL-protein standards and label-free estimated concentrations. The models with the highest accuracy were used to estimate proteome-wide label-free concentrations for all proteins from their DIA MS intensities (1,043 proteins on average within each sample; minimal two peptides per protein; see above): summing the five most intense fragment ion intensities of the most or three of the most intense peptides per protein for CO or high- $H_2$  CO cultures, respectively; summing the five most intense fragment ion intensities of all quantified peptides of the protein divided by the number of theoretically observable peptides (i.e., iBAQ<sup>70</sup>) for syngas cultures.

**Expected protein complex stoichiometries.** Equimolar stoichiometries for the HytA-E/FdhA and MetFV protein complexes were expected based on SDS gel staining experiments in *C. autoethanogenum*<sup>57</sup> and the acetogen *Moorella thermoacetica*<sup>71</sup>, respectively. Expected stoichiometries for other protein complexes in Fig. 3 and Supplementary Fig. 3 were based on measured stoichiometries in *E. coli* K-12 (Complex Portal; [www.ebi.ac.uk/complexportal](http://www.ebi.ac.uk/complexportal)) and significant homology between complex member proteins in *C. autoethanogenum* and *E. coli*. All depicted *C. autoethanogenum* protein complex members had NCBI protein-protein BLAST E-values <10<sup>-16</sup> and scores >73 against respective *E. coli* K-12 proteins using non-redundant protein sequences.

**Generation of proteomaps.** The distribution of proteome-wide protein concentrations among functional gene classifications was visualised using proteomaps<sup>27</sup>. For this, the NCBI annotation of sequence NC\_022592.1<sup>67</sup> was used as the annotation genome for *C. autoethanogenum*, with CAETHG\_RS07860 removed and replaced with the carbon monoxide dehydrogenase genes named CAETHG\_RS07861 and CAETHG\_RS07862 with initial IDs of CAETHG\_1620 and 1621, respectively. Functional categories were assigned to protein sequences with KO IDs<sup>51</sup> using the KEGG annotation tool BlastKOALA<sup>72</sup>. Since proteomaps require a tree-like hierarchy, proteins that were automatically assigned to multiple functional categories were manually assigned to one bottom-level category (Level 3 in Supplementary Table 4) based on their principal task. We also created functional categories “C1 fixation/Wood-Ljungdahl Pathway” (Level 2/3), “Acetate & ethanol synthesis” (Level 3), “Energy conservation” (Level 3), “Hydrogenases” (Level 3) and manually assigned key acetogen proteins to these categories to reflect more accurately functional categories for an acetogen. Proteins without designated KO IDs were manually assigned to latterly created categories or grouped under “Not Included in Pathway or Brite” (Level 1) with Level 2 and 3 as “No KO ID”. If BlastKOALA assigned multiple genes the same proposed gene/protein name, unique numbers were added to names (e.g., pfkA, pfkA2). The final “treemap” defining the hierarchy for our proteomaps is in Supplementary Table 4.

**Calculation of apparent in vivo catalytic rates of enzymes ( $k_{app}$ ).** We calculated apparent in vivo catalytic rates of enzymes, denoted as  $k_{app}$  ( $s^{-1}$ )<sup>26</sup>, as the ratio of specific flux rate (mmol/gDCW/h) determined before<sup>18</sup> and protein concentration (nmol/gDCW) quantified here for the same *C. autoethanogenum* CO, syngas, and high-H<sub>2</sub> CO cultures. Gene-protein-reaction (GPR) data of the genome-scale metabolic model iCLAU786<sup>18</sup> were manually curated to reflect most recent knowledge and were used to link metabolic fluxes with catalysing enzymes. For reactions with multiple assigned enzymes (i.e., isoenzymes), the enzyme with the highest average ranking of its concentration across the three cultures (Supplementary Table 3) was assumed to solely catalyse the flux. For enzyme complexes, average of quantified subunit concentrations was used (standard deviation estimated using error propagation). For the HytA-E hydrogenase, its measured protein concentration was split between reactions “rxn08518\_c0” (direct CO<sub>2</sub> reduction with H<sub>2</sub> in complex with FdhA) and “leq000001\_c0” (H<sub>2</sub> oxidation solely by HytA-E) proportionally to flux for syngas and high-H<sub>2</sub> CO cultures. The resulting enzymes or enzyme complexes catalysing specific fluxes are shown in Supplementary Table 5. Finally, we assumed each protein chain being catalytically active and only calculated  $k_{app}$  values for metabolic reactions with a non-zero flux in at least one condition, specific flux rate > 0.1% of CO fixation flux in at least one condition, and label-free data with measured concentration for the associated enzyme(s) in all conditions (Supplementary Table 5). Membrane proteins were excluded from  $k_{app}$  calculations to avoid bias from potentially incomplete protein extraction. This produced  $k_{app}$  values for 13 enzymes/complexes using anchor protein concentrations and for 48 enzymes/complexes using label-free protein concentrations (Supplementary Table 5).

**Determination of regulation level of metabolic fluxes.** We used published flux and relative proteomics data<sup>18</sup> of the same cultures studied here to determine whether fluxes are regulated by changing enzyme concentrations or their catalytic rates by considering metabolic fluxes with non-zero specific flux rates in at least two conditions of CO, syngas, or high-H<sub>2</sub> CO cultures. The same manually curated GPRs and criteria for isoenzymes and protein complexes as described above for  $k_{app}$  calculation were used to determine flux-enzyme pairs (Supplementary Table 6). We first used a two-tailed two-sample equal variance Student’s t-test with FDR correction<sup>73</sup> to determine fluxes with

significant changes between any two conditions ( $q\text{-value} < 0.05$ ). We then used the Student's left-tailed t-distribution with FDR to determine if the significant flux change for every flux was significantly different from the change in respective enzyme expression between the same conditions (Supplementary Table 6). Flux with a  $q\text{-value} < 0.05$  for the latter test was considered to be regulated at post-translational level (e.g., by changing enzyme catalytic rate).

## DATA AVAILABILITY

All data generated or analysed during this study are in the main text, supplementary information files, or public databases. Raw PRM MS data have been deposited to Panorama at [https://panoramaweb.org/Valgepea\\_Cauto\\_PRM.url](https://panoramaweb.org/Valgepea_Cauto_PRM.url) (private reviewer account details: username: panorama+reviewer27@proteinms.net; password: hSAwNUAL) with a ProteomeXchange Consortium (<http://proteomecentral.proteomexchange.org>) dataset identifier PXD025760. Raw DIA MS data have been deposited to the ProteomeXchange Consortium (<http://proteomecentral.proteomexchange.org>) via the PRIDE partner repository<sup>69</sup> with the dataset identifier PXD025732 (private reviewer account details: username: reviewer\_pxd025732@ebi.ac.uk; password: QdkJHXy0). PRM MS data with all Skyline processing settings can be viewed and downloaded from Panorama at [https://panoramaweb.org/Valgepea\\_Cauto\\_PRM.url](https://panoramaweb.org/Valgepea_Cauto_PRM.url). DIA MS data with all Skyline processing settings can be viewed and downloaded from Panorama at [https://panoramaweb.org/Valgepea\\_Cauto\\_Anchors.url](https://panoramaweb.org/Valgepea_Cauto_Anchors.url) for anchor protein absolute quantification and at [https://panoramaweb.org/Valgepea\\_Cauto\\_LF.url](https://panoramaweb.org/Valgepea_Cauto_LF.url) for estimating proteome-wide absolute protein concentrations (private reviewer account details: username: panorama+reviewer27@proteinms.net; password: hSAwNUAL). Any other relevant data are available from the corresponding author on reasonable request.

## REFERENCES

1. Liew, F. et al. Gas Fermentation—A flexible platform for commercial scale production of low-carbon-fuels and chemicals from waste and renewable feedstocks. *Front. Microbiol.* **7**, 694 (2016).



- 666 2. Claassens, N. J., Sousa, D. Z., dos Santos, V. A. P. M., de Vos, W. M. & van der Oost, J.  
667 Harnessing the power of microbial autotrophy. *Nat. Rev. Microbiol.* **14**, 692–706 (2016).
- 668 3. Fackler, N. et al. Stepping on the gas to a circular economy: accelerating development of carbon-  
669 negative chemical production from gas fermentation. *Annu. Rev. Chem. Biomol. Eng.* **12**, 1 (2021).
- 670 4. Ragsdale, S. W. & Pierce, E. Acetogenesis and the Wood–Ljungdahl pathway of CO<sub>2</sub> fixation.  
671 *Biochim. Biophys. Acta* **1784**, 1873–1898 (2008).
- 672 5. Wood, H. G. Life with CO or CO<sub>2</sub> and H<sub>2</sub> as a source of carbon and energy. *FASEB J.* **5**, 156–163  
673 (1991).
- 674 6. Drake, H. L., Küsel, K. & Matthies, C. In *The Prokaryotes* Ch. Acetogenic Prokaryotes. 354–420  
675 (2006).
- 676 7. Fuchs, G. Alternative pathways of carbon dioxide fixation: insights into the early evolution of life?  
677 *Ann. Rev. Microbiol.* **65**, 631–658 (2011).
- 678 8. Fast, A. G. & Papoutsakis, E. T. Stoichiometric and energetic analyses of non-photosynthetic CO<sub>2</sub>-  
679 fixation pathways to support synthetic biology strategies for production of fuels and chemicals. *Curr.*  
680 *Opin. Chem. Eng.* **1**, 380–395 (2012).
- 681 9. Cotton, C. A., Edlich-Muth, C. & Bar-Even, A. Reinforcing carbon fixation: CO<sub>2</sub> reduction  
682 replacing and supporting carboxylation. *Curr. Opin. Biotechnol.* **49**, 49–56 (2018).
- 683 10. Köpke, M. & Simpson, S. D. Pollution to products: recycling of ‘above ground’ carbon by gas  
684 fermentation. *Curr. Opin. Biotechnol.* **65**, 180–189 (2020).
- 685 11. Russell, M. J. & Martin, W. The rocky roots of the acetyl-CoA pathway. *Trends Biochem. Sci.* **29**,  
686 358–363 (2004).
- 687 12. Weiss, M. C. et al. The physiology and habitat of the last universal common ancestor. *Nat.*  
688 *Microbiol.* **1**, 16116 (2016).
- 689 13. Varma, S. J., Muchowska, K. B., Chatelain, P. & Moran, J. Native iron reduces CO<sub>2</sub> to  
690 intermediates and end-products of the acetyl-CoA pathway. *Nat. Ecol. Evol.* **2**, 1019–1024 (2018).
- 691 14. Ljungdahl, L. G. A life with acetogens, thermophiles, and cellulolytic anaerobes. *Annu. Rev.*  
692 *Microbiol.* **63**, 1–25 (2009).

693 15. Ragsdale, S. W. Enzymology of the Wood-Ljungdahl pathway of acetogenesis. *Ann. N. Y. Acad.*  
694 *Sci.* **1125**, 129–136 (2008).

695 16. Schuchmann, K. & Müller, V. Autotrophy at the thermodynamic limit of life: a model for energy  
696 conservation in acetogenic bacteria. *Nat. Rev. Microbiol.* **12**, 809–821 (2014).

697 17. Molitor, B., Marcellin, E. & Angenent, L. T. Overcoming the energetic limitations of syngas  
698 fermentation. *Curr. Opin. Chem. Biol.* **41**, 84–92 (2017).

699 18. Valgepea, K. et al. H<sub>2</sub> drives metabolic rearrangements in gas-fermenting *Clostridium*  
700 *autoethanogenum*. *Biotechnol. Biofuels* **11**, 55 (2018).

701 19. Valgepea, K. et al. Maintenance of ATP homeostasis triggers metabolic shifts in gas-fermenting  
702 acetogens. *Cell Syst.* **4**, 505–515.e5 (2017).

703 20. Mahamkali, V. et al. Redox controls metabolic robustness in the gas-fermenting acetogen  
704 *Clostridium autoethanogenum*. *Proc. Natl. Acad. Sci. U. S. A.* **117**, 13168–13175 (2020).

705 21. Richter, H. et al. Ethanol production in syngas-fermenting *Clostridium ljungdahlii* is controlled by  
706 thermodynamics rather than by enzyme expression. *Energy Environ. Sci.* **9**, 2392–2399 (2016).

707 22. de Souza Pinto Lemgruber, R. et al. A TetR-family protein (CAETHG\_0459) activates  
708 transcription from a new promoter motif associated with essential genes for autotrophic growth in  
709 acetogens. *Front. Microbiol.* **10**, 2549 (2019).

710 23. Song, Y. et al. Determination of the genome and primary transcriptome of syngas fermenting  
711 *Eubacterium limosum* ATCC 8486. *Sci. Rep.* **7**, 13694 (2017).

712 24. Al-Bassam, M. M. et al. Optimisation of carbon and energy utilisation through differential  
713 translational efficiency. *Nat. Commun.* **9**, 4474 (2018).

714 25. Song, Y. et al. Genome-scale analysis of syngas fermenting acetogenic bacteria reveals the  
715 translational regulation for its autotrophic growth. *BMC Genomics* **19**, 837 (2018).

716 26. Valgepea, K., Adamberg, K., Seiman, A. & Vilu, R. *Escherichia coli* achieves faster growth by  
717 increasing catalytic and translation rates of proteins. *Mol. Biosyst.* **9**, 2344–2358 (2013).

718 27. Liebermeister, W. et al. Visual account of protein investment in cellular functions. *Proc. Natl.*  
719 *Acad. Sci. U. S. A.* **111**, 8488–8493 (2014).

720 28. Schubert, O. T. et al. Absolute proteome composition and dynamics during dormancy and  
721 resuscitation of *Mycobacterium tuberculosis*. *Cell Host Microbe* **18**, 96–108 (2015).

722 29. Davidi, D. et al. Global characterisation of in vivo enzyme catalytic rates and their correspondence  
723 to in vitro kcat measurements. *Proc. Natl. Acad. Sci. U. S. A.* **104**, 20167–20172 (2016).

724 30. Calderón-Celis, F., Encinar, J. R. & Sanz-Medel, A. Standardization approaches in absolute  
725 quantitative proteomics with mass spectrometry. *Mass Spectrom. Rev.* **37**, 715–737 (2017).

726 31. Maaß, S. & Becher, D. Methods and applications of absolute protein quantification in microbial  
727 systems. *J. Proteomics* **136**, 222–233 (2016).

728 32. Ludwig, C., Claassen, M., Schmidt, A. & Aebersold, R. Estimation of absolute protein quantities  
729 of unlabeled samples by selected reaction monitoring mass spectrometry. *Mol. Cell. Proteomics* **11**,  
730 M111.013987 (2012).

731 33. Blein-Nicolas, M. & Zivy, M. Thousand and one ways to quantify and compare protein  
732 abundances in label-free bottom-up proteomics. *Biochim. Biophys. Acta* **1864**, 883–895 (2016).

733 34. Ahrné, E., Molzahn, L., Glatter, T. & Schmidt, A. Critical assessment of proteome-wide label-free  
734 absolute abundance estimation strategies. *Proteomics* **17**, 2567–2578 (2013).

735 35. Schmidt, A. et al. Absolute quantification of microbial proteomes at different states by directed  
736 mass spectrometry. *Mol. Syst. Biol.* **7**, 510 (2011).

737 36. Maier, T. et al. Quantification of mRNA and protein and integration with protein turnover in a  
738 bacterium. *Mol. Syst. Biol.* **7**, 511 (2011).

739 37. Marguerat, S. et al. Quantitative analysis of fission yeast transcriptomes and proteomes in  
740 proliferating and quiescent cells. *Cell* **151**, 671–683 (2012).

741 38. Zeiler, M., Moser, M. & Mann, M. Copy number analysis of the murine platelet proteome  
742 spanning the complete abundance range. *Mol. Cell. Proteomics* **13**, 3435–3445 (2014).

743 39. Schmidt, A. et al. The quantitative and condition-dependent *Escherichia coli* proteome. *Nat.*  
744 *Biotechnol.* **34**, 104–110 (2015).

745 40. Beck, M. et al. The quantitative proteome of a human cell line. *Mol. Syst. Biol.* **7**, 549 (2011).

746 41. Malmström, J. et al. Proteome-wide cellular protein concentrations of the human pathogen  
747 *Leptospira interrogans*. *Nature* **460**, 762–765 (2009).

748 42. Brun, V. et al. Isotope-labeled protein standards: toward absolute quantitative proteomics. *Mol.*  
749 *Cell. Proteomics* **6**, 2139–2149 (2007).

750 43. Shuford, C. M. et al. Absolute protein quantification by mass spectrometry: not as simple as  
751 advertised. *Anal. Chem.* **89**, 7406–7415 (2017).

752 44. Takemori, N. et al. High-throughput synthesis of stable isotope-labeled transmembrane proteins  
753 for targeted transmembrane proteomics using a wheat germ cell-free protein synthesis system. *Mol.*  
754 *Biosyst.* **11**, 361–365 (2015).

755 45. Takemori, N. et al. MEERCAT: multiplexed efficient cell free expression of recombinant  
756 QconCATs for large scale absolute proteome quantification. *Mol. Cell. Proteomics* **16**, 2169–2183  
757 (2017).

758 46. Gillet, L. C. et al. Targeted data extraction of the MS/MS spectra generated by data-independent  
759 acquisition: a new concept for consistent and accurate proteome analysis. *Mol. Cell. Proteomics* **11**,  
760 O111.016717 (2012).

761 47. Reiter, L. et al. mProphet: automated data processing and statistical validation for large-scale  
762 SRM experiments. *Nat. Methods* **8**, 430–435 (2011).

763 48. MacLean, B. et al. Skyline: an open source document editor for creating and analysing targeted  
764 proteomics experiments. *Bioinformatics* **26**, 966–968 (2010).

765 49. Rosenberger, G., Ludwig, C., Röst, H. L., Aebersold, R. & Malmström, L. aLFQ: an R-package  
766 for estimating absolute protein quantities from label-free LC-MS/MS proteomics data. *Bioinformatics*  
767 **30**, 2511–2513 (2014).

768 50. Lahtvee, P.-J. et al. Absolute quantification of protein and mRNA abundances demonstrate  
769 variability in gene-specific translation efficiency in yeast. *Cell Syst.* **4**, 495–504.e5 (2017).

770 51. Kanehisa, M., Sato, Y., Kawashima, M., Furumichi, M. & Tanabe, M. KEGG as a reference  
771 resource for gene and protein annotation. *Nucleic Acids Res.* **44**, D457–D462 (2016).

772 52. Peebo, K. et al. Proteome reallocation in *Escherichia coli* with increasing specific growth rate.  
773 *Mol. Biosyst.* **11**, 1184–1193 (2015).

774 53. Lemaire, O. N. & Wagner, T. Gas channel rerouting in a primordial enzyme: Structural insights of  
775 the carbon-monoxide dehydrogenase/acetyl-CoA synthase complex from the acetogen *Clostridium*  
776 *autoethanogenum*. *Biochim. Biophys. Acta* **1862**, 148330 (2021).

777 54. Liew, F. et al. Insights into CO<sub>2</sub> fixation pathway of *Clostridium autoethanogenum* by targeted  
778 mutagenesis. *mBio* **7**, e00427-16 (2016).

779 55. Liew, F. et al. Metabolic engineering of *Clostridium autoethanogenum* for selective alcohol  
780 production. *Metab. Eng.* **40**, 104–114 (2017).

781 56. Marcellin, E. et al. Low carbon fuels and commodity chemicals from waste gases – systematic  
782 approach to understand energy metabolism in a model acetogen. *Green Chem.* **18**, 3020–3028 (2016).

783 57. Wang, S. et al. NADP-specific electron-bifurcating [FeFe]-hydrogenase in a functional complex  
784 with formate dehydrogenase in *Clostridium autoethanogenum* grown on CO. *J. Bacteriol.* **195**, 4373–  
785 4386 (2013).

786 58. Mock, J. et al. Energy conservation associated with ethanol formation from H<sub>2</sub> and CO<sub>2</sub> in  
787 *Clostridium autoethanogenum* involving electron bifurcation. *J. Bacteriol.* **197**, 2965–2980 (2015).

788 59. Maia, L. B., Fonseca, L., Moura, I. & Moura, J. J. G. Reduction of carbon dioxide by a  
789 molybdenum-containing formate dehydrogenase: a kinetic and mechanistic study. *J. Am. Chem. Soc.*  
790 **138**, 8834–8846 (2016).

791 60. Schuchmann, K. & Müller, V. Direct and reversible hydrogenation of CO<sub>2</sub> to formate by a  
792 bacterial carbon dioxide reductase. *Science* **342**, 1382–1385 (2013).

793 61. Bar-Even, A. et al. The moderately efficient enzyme: evolutionary and physicochemical trends  
794 shaping enzyme parameters. *Biochemistry* **50**, 4402–4410 (2011).

795 62. Song, Y. et al. Functional cooperation of the glycine synthase-reductase and Wood–Ljungdahl  
796 pathways for autotrophic growth of *Clostridium drakei*. *Proc. Natl. Acad. Sci. U. S. A.* **117**, 7516–  
797 7523 (2020).

798 63. Tan, Y., Liu, J., Liu, Z. & Li, F. Characterization of two novel butanol dehydrogenases involved  
799 in butanol degradation in syngas-utilising bacterium *Clostridium ljungdahlii* DSM 13528. *J. Basic*  
800 *Microbiol.* **54**, 996–1004 (2014).

64. Davidi, D. & Milo, R. Lessons on enzyme kinetics from quantitative proteomics. *Curr. Opin. Biotechnol.* **46**, 81–89 (2017).
65. Nilsson, A., Nielsen, J. & Palsson, B. Ø. Metabolic models of protein allocation call for the kinetome. *Cell Syst.* **5**, 538–541 (2017).
66. Heckmann, D. et al. Machine learning applied to enzyme turnover numbers reveals protein structural correlates and improves metabolic models. *Nat. Commun.* **9**, 5252 (2018).
67. Brown, S. D. et al. Comparison of single-molecule sequencing and hybrid approaches for finishing the genome of *Clostridium autoethanogenum* and analysis of CRISPR systems in industrial relevant *Clostridia*. *Biotechnol. Biofuels* **7**, 40 (2014).
68. Wiśniewski, J. R., Zougman, A., Nagaraj, N. & Mann, M. Universal sample preparation method for proteome analysis. *Nat. Methods* **6**, 359–362 (2009).
69. Vizcaíno, J. A. et al. The Proteomics Identifications (PRIDE) database and associated tools: status in 2013. *Nucleic Acids Res.* **41**, D1063–D1069 (2012).
70. Schwanhäusser, B. et al. Global quantification of mammalian gene expression control. *Nature* **473**, 337–342 (2011).
71. Mock, J., Wang, S., Huang, H., Kahnt, J. & Thauer, R. K. Evidence for a hexaheteromeric methylenetetrahydrofolate reductase in *Moorella thermoacetica*. *J. Bacteriol.* **196**, 3303–3314 (2014).
72. Kanehisa, M., Sato, Y. & Morishima, K. BlastKOALA and GhostKOALA: KEGG tools for functional characterisation of genome and metagenome sequences. *J. Mol. Biol.* **428**, 726–731 (2016).
73. Benjamini, Y. & Hochberg, Y. Controlling the false discovery rate: a practical and powerful approach to multiple testing. *J. R. Stat. Soc. B* **57**, 289–300 (1995).

## ACKNOWLEDGEMENTS

We thank Jörg Bernhardt for help with proteomaps, Tim McCubbin for scripts, Andrus Seiman for statistics, and Olivier Lemaire for valuable discussions. This work was funded by the Australian Research Council (ARC LP140100213) in collaboration with LanzaTech. We thank the following investors in LanzaTech's technology: Sir Stephen Tindall, Khosla Ventures, Qiming Venture Partners, Softbank China, the Malaysian Life Sciences Capital Fund, Mitsui, Primetals, CICC Growth Capital

Fund I, L.P. and the New Zealand Superannuation Fund. The research utilised equipment and support provided by the Queensland node of Metabolomics Australia, an initiative of the Australian Government being conducted as part of the NCRIS National Research Infrastructure for Australia. There was no funding support from the European Union for the experimental part of the study. However, K.V. acknowledges support also from the European Union's Horizon 2020 research and innovation programme under grant agreement N810755.

## AUTHOR CONTRIBUTIONS

Conceptualisation, K.V., C.L., M.K., L.K.N., and E.M.; Methodology, K.V., G.T., N.T., A.T., C.L., A.P.M., R.T., and E.M.; Formal Analysis, K.V., and G.T.; Investigation, K.V., G.T., N.T., and A.T.; Resources, N.T., A.T., M.K., S.D.S., L.K.N., and E.M.; Writing – Original Draft, K.V. and E.M.; Writing – Review & Editing, K.V., N.T., C.L., A.P.M., R.T., M.K., L.K.N., and E.M.; Supervision, M.K., L.K.N., and E.M.; Project Administration, R.T., and E.M.; Funding Acquisition, M.K., S.D.S., L.K.N., and E.M..

## COMPETING INTERESTS

LanzaTech has interest in commercial gas fermentation with *C. autoethanogenum*. A.P.M, R.T., M.K., and S.D.S. are employees of LanzaTech.

## MATERIALS & CORRESPONDENCE

Correspondence and request for materials should be addressed to E.M.

## FIGURE LEGENDS

**Fig. 1 Absolute proteome quantification framework in *C. autoethanogenum*.** Absolute proteome quantification in light (no stable-isotope labelled [SIL] substrates) autotrophic *C. autoethanogenum* chemostat cultures was built on using 19 synthetic heavy SIL-protein spike-in standards and data-independent acquisition (DIA) mass spectrometry (MS) analysis. Culture samples with SIL-protein spike-ins and samples for DIA spectral library were analysed by DIA MS. Subsequent stringent data

analysis allowed to quantify intracellular concentrations for 16 key *C. autoethanogenum* proteins using light-to-heavy ratios between endogenous and spike-in DIA MS intensities. These 16 key proteins were further used as anchor proteins for label-free estimation of ~1,043 protein concentrations through establishing a linear correlation between protein concentrations and their measured MS intensities. Some parts created with BioRender.com.

**Fig. 2 Label-free estimation of proteome-wide protein concentrations.** **a** Correlation of peptide mass spectrometry (MS) feature intensities between biological replicate cultures of the three gas mixtures. **b** Linear correlation between anchor protein concentrations and their measured MS intensities for one syngas culture. *gDCW*, gram of dry cell weight, *aLFQ*, absolute label-free quantification. **c** Errors of different label-free quantification models for the linear fit between anchor protein concentrations and their measured MS intensities determined by bootstrapping using the aLFQ R package<sup>49</sup> for one syngas culture. *CV-MFE*, cross-validated mean fold-error. **d** Label-free quantification error of optimal model (orange) and total proteome mass (blue) across samples. Error bars denote 95% CI.

**Fig. 3 Strong correlation between protein concentrations and expected stoichiometries for equimolar protein complexes.** Grey dotted lines denote the average 1.5-fold cross-validated mean fold-error (CV-MFE) of label-free protein concentrations. Label-free protein concentrations are plotted, except for the HytA-FdhA complex, which was quantified using stable-isotope labelled protein spike-ins. Data points of the same colour represent gas mixtures. See Methods for details on expected protein complex stoichiometries. See Supplementary Table S3 for gene/protein ID, proposed name, description, and label-free data. See Table 1 for HytA-FdhA data. *gDCW*, gram of dry cell weight.

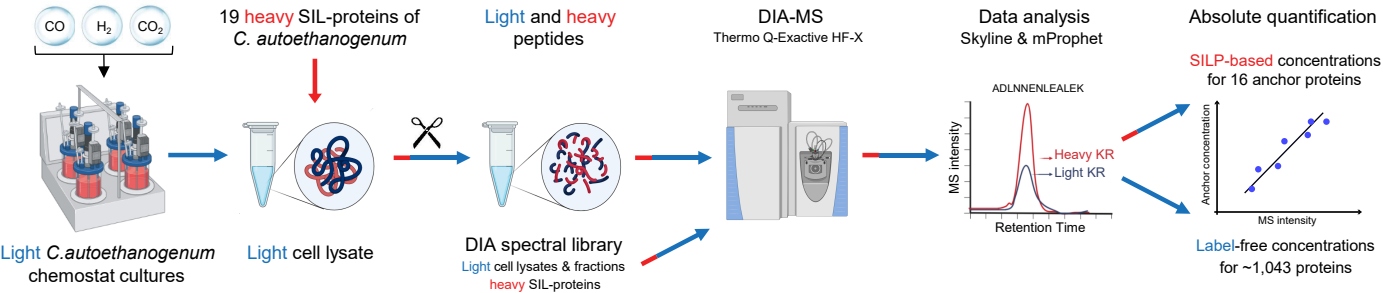
**Fig. 4 Proteomaps uncover global proteome allocation.** Left proteomap shows proteome allocation amongst functional gene classification categories (KEGG Orthology identifiers [KO IDs]<sup>51</sup>) at level two of the four-level “treemap” hierarchy structure (Supplementary Table S4). Right proteomap

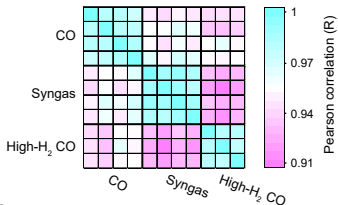
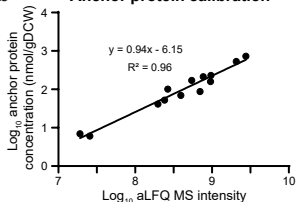
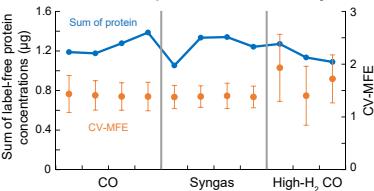


shows proteome allocation at the level of single proteins (level four of “treemap”). See Supplementary Fig. S5 for proteomaps of levels one and three of “treemap”. Area of the tile is proportional to protein concentration. Colours denote level one categories of “treemap”. Proteomaps visualise average concentrations of syngas cultures while category percentages are average of three gas mixtures (shown for categories with a fraction >5%). See Supplementary Table S3 for gene/protein ID, proposed name, description, and label-free protein concentrations.

**Fig. 5 Quantitative systems-level view of acetogen central metabolism.** Enzyme concentrations (nmol/gDCW), apparent in vivo catalytic rates of enzymes ( $k_{app}$ ;  $s^{-1}$ ), and metabolic flux rates (mmol/gDCW/h) are shown for *C. autoethanogenum* steady-state chemostat cultures grown on three gas mixtures. See dashed inset for bar chart and heatmap details. Enzyme concentration and  $k_{app}$  data are average of biological replicates. Proteins forming a complex are highlighted with non-black borders (FdhA forms a complex with HytA–E for direct  $CO_2$  reduction with  $H_2$ ; CooS1 is expected to form a complex with CooS1a and b as they are encoded from the same operon). For reactions with isoenzymes,  $k_{app}$  is for the enzyme with the highest concentration ranking (top location on enzyme heatmap), see Methods for details. Flux data from ref.<sup>18</sup> are average of biological replicates and error bars denote standard deviation. Arrows show direction of calculated fluxes; red arrow denotes uptake or secretion. Gene/protein IDs right of enzyme concentration heatmaps are preceded with CAETHG\_RS and red font denotes concentrations determined using stable-isotope labelled (SIL) protein spike-in standards (i.e., anchor proteins). Asterisk denotes data for redox-consuming  $CO_2$  reduction to formate solely by FdhA without the use of  $H_2$  during growth on  $CO$ . <sup>a</sup>Bifunctional acetaldehyde/alcohol dehydrogenase (acetyl-CoA→ethanol); <sup>b</sup>Flux into PEP from OAA and pyruvate is merged and  $k_{app}$  is for PEPCK. See Supplementary Table S3 for gene/protein ID, proposed name, description, and label-free protein concentrations. See Table 1 for anchor protein concentrations. See Supplementary Table S5 for  $k_{app}$  and flux data. See ref.<sup>18</sup> for cofactors of reactions and metabolite abbreviations. gDCW, gram of dry cell weight, *NQ*, not quantified.

912 **Fig. 6 Regulatory principles of apparent in vivo catalytic rates of enzymes ( $k_{app}$ ) and metabolic**  
 913 **flux throughput. a** Enzymes catalysing higher metabolic flux rates have both higher concentrations  
 914 and higher  $k_{app}$ s. Yellow and blue denote high and low values, respectively. Kendall's  $\tau$  correlations  
 915 with significance p-values between respective pairs are shown below heatmap. See Supplementary  
 916 Table S5 for flux rate, enzyme concentration, and  $k_{app}$  data, and for description of reaction names  
 917 (Rxn name) and gene-protein-reaction (GPR) associations. **b** Control of metabolic flux throughput  
 918 through  $k_{app}$  changes for high flux pathways. See also Fig. 5. *gDCW*, gram of dry cell weight.



**a****Correlation of peptide features****b****Anchor protein calibration****d****Label-free quantification accuracy****c****aLFQ optimal model determination**