

1

2 **Supplementary Information for**
3 **Boundary conditions for early life converge to an organo-sulfur metabolism**

4 **Joshua E. Goldford, Hyman Hartman, Robert Marsland III and Daniel Segrè**

5 **Daniel Segrè**
6 **E-mail: dsegre@bu.edu**

7 **This PDF file includes:**

- 8 Supplementary text
- 9 Figs. S1 to S4
- 10 Caption for Database S1
- 11 References for SI reference citations

12 **Other supplementary materials for this manuscript include the following:**

- 13 Database S1

14 Materials and methods

15 **Reconstruction of biosphere-level metabolic network.** Biosphere-level metabolism was reconstructed from the KEGG database
16 (1) according to protocol described previously (2). We modified the network in several ways to model primitive thioester-based
17 metabolic network without nitrogen or phosphate. First, to simulate the availability of thiols capable of forming thioesters, we
18 included Coenzyme A, Acyl-Carrier Protein and Glutathione into the seed set. However, to enforce the constraint that these
19 metabolites could only be used in reactions as coenzymes (and not products or substrates), we prevented the degradation by
20 removing KEGG reactions R10747, R02973 and R02972.

21 We next assigned standard molar free energies to reactions using eQuilibrator at a predefined pH (3). Next we substituted
22 NAD, NADP and FAD-coupled reactions with an arbitrary redox couple. For example, if the redox reaction $X_{ox} + \text{NADH} \rightarrow$
23 $X_{red} + \text{NAD}^+$ was swapped with electron donor with a redox potential of E_0^+ mV, we would use the following formula to
24 adjust the standard molar free energy for the new reaction r' :

$$25 \quad \Delta_{r'}G'^{\circ} = \Delta_rG'^{\circ} + nF(E_0^+ - E_0)$$

26 where n are the number of electrons transferred in reaction r and $F = 96.485$ kJ/V. Note that if we assumed that the
27 electron donor/acceptor substitute was a two electron donor/acceptor, we did not change the stoichiometry in the reaction
28 equation. However, in the case where the electron donor/acceptor substitute was a single electron donor/acceptor, we change
29 the stoichiometric coefficients to $s_{cj} = 2$ for all reactions j , where c represents metabolites NAD(H), NADP(H) and FAD(H₂).
30 For this work, we systematically varied the reduction potential E_0^+ and stoichiometry of the primitive redox coenzyme.

31 **Thermodynamically-constrained network expansion.** We performed network expansion using thermodynamic constraints in
32 a different way than performed previously (2) Previously, we removed reactions above a predefined free energy threshold
33 of $\tau = 30$ kJ/mol (2). For this work, we computed the lowest reaction free energy possible using estimates for upper and
34 lower bounds on metabolite concentrations, u_i and l_i , and removed reactions with a positive reaction free energy. For a given
35 biochemical reaction at fixed temperature and pressure, Δ_rG' is defined as:

$$36 \quad \Delta_rG' = \Delta_rG'^{\circ} + RT \ln \prod_i a_i^{s_{ir}}$$

37 where the $\Delta_rG'^{\circ}$ is the free energy change of the reaction at standard molar conditions, R is the ideal gas constant, T is
38 temperature, a_i is the activity of metabolite i and s_{ir} is the stoichiometric coefficient for metabolite i in reaction r . We fixed
39 a_i for each reactions according to the following rules:

$$40 \quad a_i = \begin{cases} u_i, & \text{if } s_{ir} < 0, \\ l_i, & \text{if } s_{ir} > 0. \end{cases} \quad [1]$$

41 We then removed reactions with a $\Delta_rG' \geq 0$. For all simulations we assumed that $u_i = 10^{-1}$ M and $l_i = 10^{-6}$ M. Note that
42 because we model each reaction independently, metabolite concentrations could be inconsistent. For instance, if metabolite i is
43 the substrate for reaction a and a product for reaction b , then $x_i = u_i$ for reaction a and $x_i = l_i$ for reaction b .

44 Using this procedure to systematically remove reactions that were considered to be thermodynamically infeasible, we
45 performed network expansion (4–6) using the computational procedure described in (2).

46 **Parameters for network expansion.** We systematically studied the size and composition of networks under precise environmental
47 conditions by varying (a) the reduction potential from the environment, (b) pH, (c) temperature, (d) the presence or absence
48 thiols, (e) the inclusion of fixed carbon into the seed set and (f) the inclusion of fixed nitrogen into the seed set. We now
49 discuss each of these parameters in more detail:

- 50 • *Reduction potential and stoichiometry.* A wide range of environmental conditions could have provided electron donors
51 at various potentials: high potential redox pairs, with strong oxidants like Fe(III), may have been present in oceans at
52 high concentrations, while strong reductants like H₂, disulfides, proto-ferredoxin, or reductive carboxylation of thioesters
53 have been produced via serpentinization or geochemical analogues of primitive metabolic pathways (7). We substituted
54 reactions coupled to NAD, NADP and FAD with a generic single or double electron donor and acceptor pair at a fixed
55 potential. To prevent unbalanced electron transfer, we removed the following transhydrogenase reactions: R10159, R01195,
56 R00112, R09520, R09748, R05705, R05706, R09662, R09750. We then created a single or double electron donor/acceptor
57 pair with a fixed reduction potential, E_0^+ , ranging from -600 to 600 mV.
- 58 • *pH* We modified the pH by setting reaction free energies at various pH's (5.0-9.0) using eQuilibrator (3) which relies on
59 the component contribution method (8).
- 60 • *Temperature.* Temperatures were assumed to have been within a range of 50-150 °C, spanning estimates of ocean seawater
61 temperature in the Archean (9), up to some alkaline hydrothermal vent systems (10).

62 • *Thiols*. In our model we provided thiols that serve as substitutes for coenzymes that form thioester bonds in extant
63 metabolic networks. To this end, we provided Coenzyme A, acyl-carrier protein and Glutathione in the seed set, but
64 removed key degradation reactions to ensure these compounds only served as coenzymes, rather than material sources,
65 during network expansion (2).

66 • *Fixed nitrogen*. To study the consequences of adding or removing a source of fixed nitrogen as a seed compounds for
67 network expansion, we either added or removed ammonia from the seed set prior to expansion.

68 In addition to parameters we varied, we kept constant two additional parameters that could be studied in future work:

69 • *Metabolite concentrations*. Metabolite concentrations were assumed to be within 1 μM - 100 mM. The upper bound
70 estimate is consistent with recent experimental data showing that key metabolites (formate, methanol, acetate and
71 pyruvate) can be produced near 100 mM (11). Although we do not have empirical evidence to suggest a reasonable lower
72 bound on metabolite concentrations in ancient metabolic networks, we assumed that 1 μM , the estimated lower bound in
73 today's cells (12), was also the lower bound in our model of ancient metabolism.

74 • *Reactions with no free energy estimate*. 53 % of the biosphere-level metabolic network reactions have no free energy
75 estimate (4851 of 9074). For all simulations presented in this paper, we assumed these reactions were blocked and did not
76 include them in the network.

77 **Generalized linear modeling of network expansion results.** To access the effects of various parameters on the outcome of
78 network expansion, we used generalized linear models to construct logistic regression classifiers to predict whether or not
79 the network expanded beyond 100 metabolites using a combination of predictors, including categorical variables encoding
80 whether or not ammonia, thiols or fixed carbon was provided in the seed set, and continuous variables encoding the reduction
81 potential, pH and temperature used in each simulation. We first define the response variable for simulation i as y_i where $y_i = 1$
82 if the simulation resulted in a network that expanded beyond 100 metabolites, and zero otherwise. For the set of simulations
83 performed in Fig. 1 in the main text, we constructed a design matrix consisting of categorical variables representing the
84 following scenarios:

- 85 1. $x_{N,i} \in \{0, 1\}$: 1 if ammonia was included in the seed set, and 0 otherwise.
- 86 2. $x_{S,i} \in \{0, 1\}$: 1 if thiols were included in the seed set, and 0 otherwise.
- 87 3. $x_{C,i} \in \{0, 1\}$: 1 if fixed carbon (acetate/formate) was included in the seed set, and 0 otherwise.
- 88 4. $x_{H,i} \in \mathbb{R}_{>0}$: A continuous variable representing the pH. Note for our simulations, we only explored acidic (pH=5), neutral
89 (pH=7) and alkaline (pH=9) regimes.
- 90 5. $x_{E,i} \in \mathbb{R}$: A continuous variable representing the reduction potential at standard molar conditions (at the specified pH
91 listed above). For our simulations, we explored a wide range of standard molar reduction potentials (from -600 mV to
92 +600 mV).
- 93 6. $x_{T,i} \in \mathbb{R}$: A continuous variable representing the temperature. For our simulations, we explored two temperatures: a
94 high temperature regime ($T = 150$ °C), and a low temperature regime ($T = 50$ °C).

95 We next constructed the following generalized linear model to model whether the network expanded beyond metabolites:

$$96 \quad \text{logit}(y_i) = \beta_0 + \beta_N x_{N,i} + \beta_S x_{S,i} + \beta_C x_{C,i} + \beta_H x_{H,i} + \beta_E x_{E,i} + \beta_T x_{T,i} \quad [2]$$

97 We fit the parameters (β) using the *fitglm.m* function in MATLAB 2015a, and a receiver operating curve (ROC) was
98 generated using the *perfcurve.m* function. For results presented in Fig. 2C in the main text, individual predictors were removed
99 in the generalized linear model presented above. To access whether the trained logistic model served as an accurate classifier,
100 we performed leave-one out cross-validation by removing individual samples from the training set and testing the accuracy of
101 the trained classifier on the removed sample. This procedure resulted in a cross-validation accuracy of 0.89.

102 **Constraint-based modeling.** We constructed a model of a autocatalytic network at steady state using a variant of constraint-
103 based modeling of cellular metabolism called thermodynamic-based metabolic flux analysis (TMFA) (13). TMFA transforms
104 the non-linear constraints induced by imposing thermodynamic consistency into mixed-integer linear constraints. In this
105 section, we first describe (a) the construction of primitive biomass composition for a model of an ancient proto-cell and (b) the
106 formulation of TMFA used in this analysis.

107 **Prebiotic biomass equation.** We constructed a simple model for the macromolecular composition of primitive proto-cells, using
 108 empirical knowledge of extant cellular life. Since our metabolic model of proto-metabolism does not included macromolecular
 109 production of nucleotides (and thus a nucleic acid based genetic system), we assume that the primary role of proto-cellular
 110 metabolism was to initially produce components for a cellular membrane and catalysts. Building off of Christian de Duve’s
 111 multimer hypothesis (14), we first propose that the biomass can be constructed using a simple two parameter model consisting
 112 of the mass fraction of lipids ϕ_L and the average length of each catalytic multimer n .

- 113 • *Lipid mass fraction.* The lipid content in modern cells is roughly 10% of the total dry mass (Bionumbers ID: 111209)
 114 (15), primarily composed of the fatty acid palmitate. For our analysis, we assume that palmitate represents the sole
 115 component of lipids. Future models could incorporate glycerol, which enables the production of glycerolipids. While
 116 phosphate is used in cellular membranes as a polar head group to produce amphiphilic molecules, primitive processes may
 117 have conjugated negatively charged organic acids (e.g. oxalate) to glycerol via a thioester-mediated synthesis mechanism
 118 to create amphiphilic lipid molecules resembling modern phospholipids. For our initial model, we simply propose that
 119 palmitate was the initial amphiphilic component of primitive membranes, where the negatively charged polar carboxylate
 120 ion was sufficient for forming a membrane, and assumed that proto-cells consisted of a lipid mass fraction of ϕ_L .
- 121 • *Catalytic multimer mass fraction.* We propose that ancient catalysts were composed of inorganic molecules (e.g. iron-sulfur
 122 clusters, metal ions, mineral surfaces) chelated with multimers of α -hydroxy-acids (see Fig. 4A in main text). For our
 123 model, we assume that the eight keto acid precursors producible from our network were the dominate monomers of
 124 ancient multimeric catalysts. We assume that the total mass fraction of these catalysts are $1 - \phi_L = \phi_C = \sum_k \phi_k$, where
 125 ϕ_k is the mass fraction of polymerized monomer k . For our analysis, we assumed that each monomers is uniformly
 126 distributed, so that $\phi_k = \text{constant}$ for all k . Additionally, since each monomer must be reduced to α -hydroxy acids,
 127 there is linear relationship between the electron demand, s_e , and the number of molecules of monomers produced. The
 128 stoichiometric equivalents of electron donors are thus:

$$129 \quad s_e = 2 \sum_k \frac{\phi_k}{M_k}$$

130 where M_k is the molar mass of monomer k .

- 131 • *Average size of catalytic multimers.* The average size of mulimeric catalysts sets the number of thioester bonds required
 132 for synthesis of catalytic multimers. For each polymer of size n , there are $n - 1$ thioester bonds required. In our model,
 133 the total number of monomers are fixed to be: $\sum_k \frac{\phi_k}{M_k}$, where M_k is the molecular weight for monomer k . Thus for a
 134 fixed monomer length n , we can compute the number polymers using the following formula:

$$135 \quad P(n) = \frac{1}{n} \sum_k \frac{\phi_k}{M_k}$$

136 The thioester demand is thus $s_t(n) = (n - 1)P(n)$, or:

$$137 \quad s_t(n) = \frac{n - 1}{n} \sum_k \frac{\phi_k}{M_k}$$

138 For our analysis we assumed a fixed polymer length of size $n = 10$ monomers.

139 Using these two parameters, we constructed the biomass equation for the proto-cellular model.

140 **Thermodynamic Metabolic Flux Analysis (TMFA).** To simulate a thermodynamically-feasible steady-state behavior of this metabolic
 141 network, we used thermodynamic metabolic flux analysis (TMFA) (13). Briefly, TMFA transforms the non-linear constraints
 142 induced by imposing thermodynamic consistency into mixed-integer linear constraints. We first converted the model into an
 143 irreversible model by modeling each reaction as both forward and backward half reactions. We then constructed the following
 144 mixed-integer linear program (MILP) to find a flux vector, v (with elements v_r for each reaction r), log-transformed metabolite
 145 concentrations ($\ln(x)$) and binary variables indicating whether a reaction is feasible (z) given a specific objective function was
 146 satisfied. The objective function used in this work was to maximize biomass yield, similar to the objectives frequently used in
 147 FBA model of microbial metabolism. Thus, the optimization problem was constructed according the following MILP:

$$\begin{aligned}
 & \underset{\ln(x), v, z, e}{\text{maximize}} && v_{\text{biomass}} \\
 & \text{subject to} && Sv = 0 \\
 & && 0 < v_r \leq z_r ub_r, \forall r \in \mathcal{R} \\
 & && z_r K - K + \Delta_r G' < 0, \forall r \in \mathcal{R} \\
 & && \Delta_r G'^o + RT \sum_i s_{ir} \ln(x_i) + \sigma_r e_r = \Delta_r G' \\
 & && \ln(10^{-6}) \leq \ln(x_i) \leq \ln(10^{-1}), \forall i \in \mathcal{M} \\
 & && -\sigma_m \leq e_r \leq \sigma_m \forall r \in \mathcal{R}
 \end{aligned} \tag{3}$$

149 where \mathcal{R} and \mathcal{M} are the sets of all reactions and metabolites, respectively. As discussed in detail elsewhere (13), the first
150 equation in the constraint set ensures that intracellular metabolite concentrations are at steady-state, and are simply mass
151 balance constraints for each metabolite. The second equation sets the bound on individual reaction fluxes, where the maximum
152 flux through reaction r is ub_r . Note that when $z_r = 0$, the flux through reaction r is constrained to 0. The third equation
153 sets ensures that $z_r = 1$ if and only if $\Delta_r G' < 0$, and $z_r = 0$ otherwise. Note that K is a large number ($K > \max_r \{\Delta_r G'\}$)
154 ensuring that this constraint is not violated with $z_r = 0$. The fourth equation is the free energy of each reaction as a function
155 of log-metabolite concentrations. Note that we also add slack variables, e_r , to account for the possible error in the estimating
156 standard molar reaction free energies for each reaction (where σ_r is the standard error for each reaction r), which are bounded
157 by a global error tolerance $\sigma_m = 2$ (set in equation 6). Lastly, equation 5 simply constrains the log-metabolite concentrations to
158 be bounded between $1\mu\text{M}$ and 100 mM . After each simulation, we performed a secondary optimization to find the minimal set
159 of reactions that achieve the optimal growth rate by minimizing the l_1 -norm of the flux distribution subject to the constraint
160 that $v_{\text{biomass}} = v_{\text{biomass}}^*$

161 Numerical simulations were performed using the COBRA toolbox (16) and the Gurobi optimizer (Version 7.0.1). All source
162 code is provided in the following github repository: http://www.github.com/jgoldford/protometabolic_modeling.

163 **Calculation of coenzyme and sequence-level features within enzymes.** To determine which reactions were associated with
164 specific coenzymes (for results presented in Fig. S1B,D) we downloaded information for each Enzyme Commission number
165 (E.C.) in the KEGG ENZYME database (<http://www.genome.jp/kegg/annotation/enzyme.html>). We downloaded each page and
166 parsed the "comment" field for each E.C. and performed a text-based search to identify coenzymes associated with each E.C.
167 number. We searched for text indicating that the enzyme mechanisms used one of the following coenzymes, cofactors and iron
168 sulfur clusters: biotin, heme, PLP, TPP, pterin, molybdopterin, flavin, Fe, Co, Ni, Cu, Mn, W, Zn, Mo, Mg, FeS, FeFe, Fe₂S₂,
169 Fe₃S₄ and Fe₄S₄, respectively. We also searched E.C. numbers indicating that the reaction mechanisms are non-enzymatic.
170 Text-based searches were pruned manually to remove mis-annotated enzyme-coenzyme relationships.

171 For results presented in S1B, we computed the fraction of reaction E.C. numbers that were associated with one of the
172 following coenzymes: Fe, Co, Ni, Cu, Mn, W, Zn, Mo, Mg, FeS, FeFe, Fe₂S₂, Fe₃S₄ and Fe₄S₄, or was marked as non-enzymatic.
173 For results presented in S1D, we computed the fraction of reaction E.C. numbers that were associated with one of the following
174 coenzymes: biotin, heme, PLP, TPP, pterin, molybdopterin, and flavin.

175 For results presented in Fig. S1E, we obtained a database of known enzyme active site residues (17). We first mapped the
176 network reactions to E.C. numbers listed in KEGG, then identified active sites corresponding to E.C. numbers within the
177 expanded network. We next computed the fraction of active site residues containing nitrogenous side-chains, derived from the
178 following amino acids: Arginine (R), Lysine (K), Glutamine (Q), Asparagine (N), Histidine (H), and Tryptophan (W).

Supporting Information Text

Network enzymes retain features of nitrogen-free catalysts. In order for the expanded networks presented in the main text to have operated in prebiotic conditions, reactions would have been catalyzed non-enzymatically by inorganic or simple organic catalysts available in prebiotic environments. Prior work has suggested that reactions in metabolic networks that proceed spontaneously or depend on enzymes with inorganic coenzymes, such as iron-sulfur or transition metal cofactors, may have operated in prebiotic conditions (11, 18–20). We identified reactions in KEGG that could proceed spontaneously or are dependent on one of several inorganic coenzymes (Methods), and defined this set of reactions as *plausibly pre-enzymatic*, or “PPE”-reactions (Fig. S1A). For each proposed prebiotic scenario that lead to expansion of at-least 100 metabolites ($n = 144$, Fig S1A), we partitioned reactions added to the network before the inclusion of ammonia into the seed set (herein called “pre-ammonia” reactions) and reactions added to the network after ammonia was added to the seed set (or “post-ammonia” reactions). We then computed the fraction pre- and post-ammonia reactions that were classified as PPE reactions, and found that pre-ammonia reactions contained a higher proportion of PPE-reactions relative to post-ammonia reactions (one-tailed Wilcoxon sign-rank test: $P < 10^{-19}$), suggesting that the pre-ammonia reactions may have been more readily catalyzed by simple inorganic catalysts in prebiotic environments. We next hypothesized that if these enzymes evolved from a thioester-driven proto-metabolism without nitrogen, then enzymes in these networks should be depleted in enzyme-bound nitrogen-containing coenzymes. We thusly computed the fraction of pre- and post-ammonia reactions dependent on enzymes containing TPP, PLP, heme, biotin, flavin, pterin, and cobalamin (Fig. S1C, Methods). We found that the proportion of pre-ammonia reactions associated with these coenzymes were significantly less than the proportion of post-ammonia reactions dependent on these coenzymes (Fig. S1D, one-tailed Wilcoxon sign-rank test: $P < 10^{-24}$), which is primarily due to the large number of PLP-dependent reactions added to the network after the inclusion of ammonia.

Since only a minority of reactions in this network were categorized as PPE, simple organic or organosulfur catalysts may have been necessary in order for this network to function in prebiotic environments. Christian de Duve suggested that thioester-based polymers may have provided the necessary catalytic components of ancient metabolism in addition to inorganic catalysts (14). In modern living systems, monomers of keto acids are converted into amino acids, which are then polymerized into polypeptides either with or without the aid of the ribosome and mRNA. If prebiotic environments were severely nitrogen limited, keto acids may have been reduced to hydroxy acids, and polymerized into polyesters using thioesters as a condensing agent. Notably, in such a scenario only the polymer backbone is altered, leaving the side chains (*R*-groups) within today’s amino acids intact. Recent work has demonstrated that polyesters may aid in the polymerization of amino acids during dry-wet cycles (21), and that the peptidyl-transferase domain on the ribosome can polymerize hydroxyacylated tRNAs to form polyesters (22, 23), suggesting that ester bond formation may have proceeded amide bond formation in living systems.

It has been proposed that enzymes retain features of early catalysts before the emergence of the genetic code and protein translation systems, and that enzyme active sites may bear resemblance to ancient catalysts. Thus, if this network represents a relic of an ancient metabolism before the biological incorporation of nitrogen, then the active sites of enzymes catalyzing reactions within the network should be depleted in amino acids with side chains containing nitrogen (Fig. S1E). To see if the catalytic residues of the enzymes in the pre-ammonia network were depleted in amino acids with nitrogenous side chains, we first obtained a database of catalytic site residues inferred from protein structures (17). After removing entries with interactions mediated by the peptide backbone, this dataset consisted of 18,721 entries, 1,304 of which were associated with active sites of enzymes in the nitrogen-free network in a representative network. For each putative prebiotic scenario resulting in an expansion with more than 100 metabolites, we computed the fraction of active site residues that contained nitrogen in enzymes associated with both pre- and post-ammonia reactions (Fig. S1E). We found that the proportion of nitrogenous catalytic residues associated with pre-ammonia reactions was significantly lower than the proportion of nitrogenous catalytic residues associated with post-ammonia reactions (Fig. S1F, Wilcoxon sign-rank test: $P < 10^{-24}$).

One potential alternative explanation for these biases in amino acid composition within the active sites of extant enzymes may be the outcome of evolutionary selection: nitrogen limitation in the environment may have favored mutations that lead to less nitrogen within these enzymes. However, evidence for selection for less nitrogen usage would manifest within the entire protein sequence, rather than just the active sites. Thus, we computed the fraction of amino acids with nitrogenous side chains across the entire coding sequences, rather than specifically the active sites, for enzymes associated with pre- and post-ammonia reactions (see Methods). We found no evidence that enzymes in the pre-ammonia network had a decreased usage of amino acids with nitrogenous side chains side chains relative to enzymes added to the network after ammonia was included in the seed set (one-tailed Wilcoxon sign rank test: $P = 1$), suggesting that the biases within the active sites are not merely a consequence evolutionary selection (see Fig. S4).

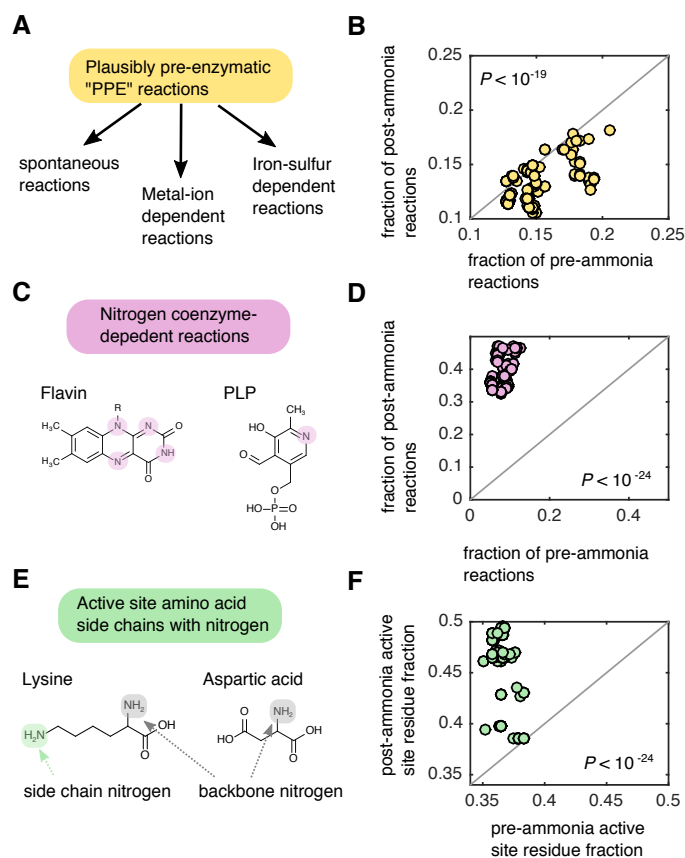


Fig. S1. Enzymes in thioester-driven protometabolism are depleted in nitrogenous compounds (A) We classified reactions in KEGG as being plausibly pre-enzymatic (PPE) reactions if they could (a) proceed spontaneously, (b) were associated with enzymes that contain at-least one iron-sulfur cluster or (c) were associated with an enzyme that relied on at-least one metal (Ni, Co, Cu, Mg, Mn, Mo, Zn, Fe, W) cofactor. (B) For all scenarios resulting in expansion of > 100 metabolites ($n = 144$, Fig. S1A) we computed the fraction of PPE-reactions amongst the pre-ammonia reactions (x axis) and post-ammonia reactions (y-axis). The frequency of PPE-reactions in the pre-ammonia reaction set was on average higher than the frequency of PPE-reactions in the post-ammonia reaction set (one-tailed Wilcoxon sign-rank test: $P < 10^{-19}$). (C) We identified KEGG reactions that were dependent on at-least one of the following nitrogen-containing coenzymes: flavin, biotin, thiamine pyrophosphate (TPP) pyridoxal phosphate (PLP), heme, pterin or cobalamin. (D) We compute the fraction of pre- and post- ammonia reactions associated with nitrogen containing coenzymes in the KEGG database, and found that a much higher proportion of post-ammonia reactions were dependent on these coenzymes relative to pre-ammonia reactions (one-tailed Wilcoxon sign-rank test: $P < 10^{-24}$). (E) We parsed the catalytic active site database (17) to find entries associated with pre and post-ammonia reactions, and compute the fraction of entries associated with amino acids with nitrogen-containing side chains (Q,N,W,H,K,R). (F) For each scenario, the fraction of active sites with nitrogen-containing amino acids was significantly higher for post-ammonia reactions relative to pre-ammonia reactions (one-tailed Wilcoxon sign-rank test: $P < 10^{-24}$).

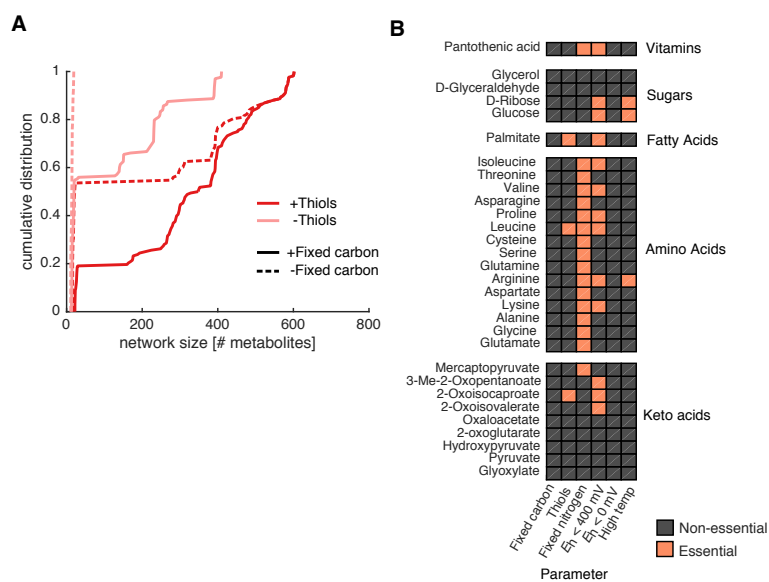


Fig. S2. Thiols are required for autotrophic expansion and fatty acid production (A) We grouped the $n = 672$ geochemical scenarios into wither a source of fixed carbon and thiols was provided in the seed set. We then plotted the empirical cumulative distributions for each group of scenarios. Notably, when thiols and fixed carbon are not supplied in the seed set, the networks are always below 100 metabolites, indicating that expansion is prohibited without either fixed carbon or thiols in the seed set. (B) We determined what geochemical parameters (x -axis) were essential for the production of important biomolecules (y -axis). For example, palmitate, a long chain fatty acid, is producible only if thiols and reductant below 400 mV is provided in the seed set.

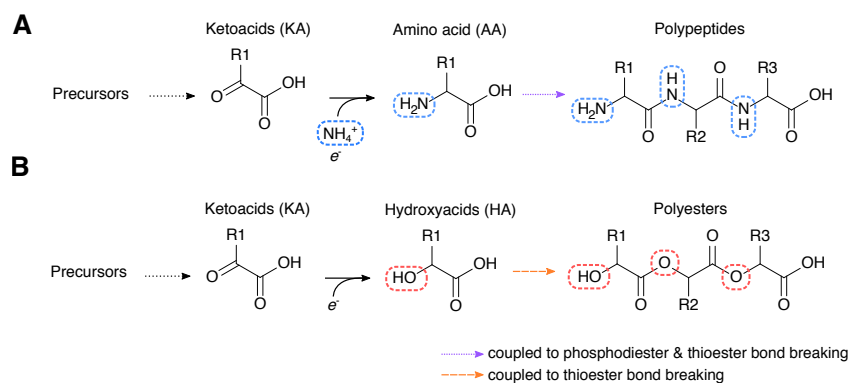


Fig. S3. Putative ancient catalysts. (A) In extant biochemistry, keto acids are converted to amino acids using transamination or reductive amination reaction mechanisms, which are then polymerized using a phosphate or thioester coupled mechanism to make polypeptides. (B) If prebiotic environments did not have a source of fixed nitrogen, then keto acids could have been reduced to α -hydroxy acids, which could then be polymerized into polyesters either with (14) or without (21) thioester bond breaking.

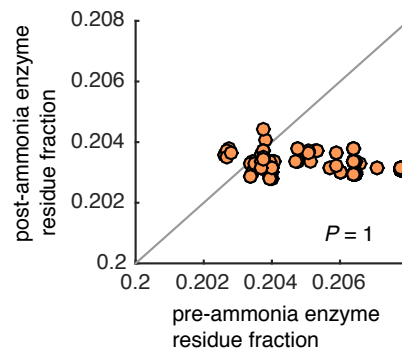


Fig. S4. Enzymes catalyzing reactions before the addition of ammonia are not depleted in nitrogen containing amino acids relative to enzymes added after ammonia To see if the amino acid biases in active sites of enzymes catalyzing reactions added to the network without ammonia (see Fig. S1E-F) is confounded due to evolutionary selection for reduced nitrogen in these enzymes, we computed the fraction of nitrogen side chains in enzymes in pre-ammonia reactions (x -axis) and in enzymes in post-ammonia reactions y -axis. We found that enzymes in the pre-ammonia networks did not have significantly less nitrogen usage compared to enzymes in post-ammonia reactions (one-tailed Wilcoxon sign-rank test: $P = 1$).

230 **Additional data table S1 (scenario_results.csv)**

231 Results of the systematic exploration of 672 prebiotic scenarios. Columns labeled with KEGG compounds IDs denote
232 whether a compound appeared after expansion (1) or not (0).

233 **References**

- 234 1. Kanehisa M, Goto S (2000) KEGG: kyoto encyclopedia of genes and genomes. *Nucleic acids research* 28(1):27–30.
235 2. Goldford JE, Hartman H, Smith TF, Segrè D (2017) Remnants of an Ancient Metabolism without Phosphate. *Cell*
236 168(6):1126–1134.e9.
237 3. Flamholz A, Noor E, Bar-Even A, Milo R (2012) EQuilibrator - The biochemical thermodynamics calculator. *Nucleic*
238 *Acids Research* 40(D1):770–775.
239 4. Ebenhöf O, Handorf T, Heinrich R (2004) Structural analysis of expanding metabolic networks. *Genome informatics.*
240 *International Conference on Genome Informatics* 15(1):35–45.
241 5. Handorf T, Ebenhöf O, Heinrich R (2005) Expanding metabolic networks: scopes of compounds, robustness, and evolution.
242 *Journal of molecular evolution* 61(4):498–512.
243 6. Raymond J, Segrè D (2006) The effect of oxygen on biochemical networks and the evolution of complex life. *Science (New*
244 *York, N. Y.)* 311(5768):1764–7.
245 7. Martin WF, Thauer RK (2017) Energy in Ancient Metabolism. 168(6):953–955.
246 8. Noor E, Haraldsdóttir HS, Milo R, Fleming RMT (2013) Consistent estimation of Gibbs energy using component
247 contributions. *PLoS computational biology* 9(7):e1003098.
248 9. Halevy I, Bachan A (2017) The geologic history of seawater pH. *Science* 355(6329):1069–1071.
249 10. Martin W, Russell MJ (2007) On the origin of biochemistry at an alkaline hydrothermal vent. *Philosophical Transactions*
250 *of the Royal Society of London B: Biological Sciences* 362(1486):1887–1926.
251 11. Varma SJ, Muchowska KB, Chatelain P, Moran J (2018) Native iron reduces CO₂ to intermediates and end-products of
252 the acetyl-CoA pathway. *Nature Ecology & Evolution* 2.
253 12. Bar-Even A, Flamholz A, Noor E, Milo R (2012) Thermodynamic constraints shape the structure of carbon fixation
254 pathways. *Biochimica et biophysica acta* 1817(9):1646–59.
255 13. Henry CS, Broadbelt LJ, Hatzimanikatis V (2007) Thermodynamics-based metabolic flux analysis. *Biophysical journal*
256 92(5):1792–1805.
257 14. de Duve C (1991) *Blueprint for a cell: the nature and origin of life.* (Neil Patterson Publishers, Carolina Biological Supply
258 Company, Burlington, N.C.), p. 275.
259 15. Milo R, Jorgensen P, Moran U, Weber G, Springer M (2010) BioNumbers—the database of key numbers in molecular and
260 cell biology. *Nucleic acids research* 38(Database issue):D750–3.
261 16. Schellenberger J, et al. (2011) Quantitative prediction of cellular metabolism with constraint-based models: the COBRA
262 Toolbox v2.0. *Nature Protocols* 6(9):1290–1307.
263 17. Ribeiro AJM, et al. (2018) Mechanism and Catalytic Site Atlas (M-CSA): a database of enzyme reaction mechanisms and
264 active sites. *Nucleic Acids Research* 46(D1):D618–D623.
265 18. Sousa FL, et al. (2013) Early bioenergetic evolution. *Philosophical transactions of the Royal Society of London. Series B,*
266 *Biological sciences* 368(1622):20130088.
267 19. Muchowska KB, et al. (2017) Metals promote sequences of the reverse Krebs cycle. *Nature Ecology and Evolution*
268 1(11):1716–1721.
269 20. Sousa FL, Preiner M, Martin WF (2018) Native metals, electron bifurcation, and CO₂ reduction in early biochemical
270 evolution. *Current opinion in microbiology* 43:77–83.
271 21. Forsythe JG, et al. (2015) Ester-Mediated Amide Bond Formation Driven by Wet-Dry Cycles: A Possible Path to
272 Polypeptides on the Prebiotic Earth. *Angewandte Chemie - International Edition* 54(34):9871–9875.
273 22. Fahnestock S, Rich A (1971) Ribosome-Catalyzed Polyester Formation. *Science* 173(3994):340–343.
274 23. Ohta A, Murakami H, Suga H (2008) Polymerization of alpha-hydroxy acids by ribosomes. *Chembiochem : a European*
275 *journal of chemical biology* 9(17):2773–2778.