

Supporting Information for

A data-driven sequencer that unveils latent “codons” in synthetic copolymers.

Yusuke Hibi^{1*}, Shiho Uesaka¹ and Naito Masanobu^{1*}

Correspondence to: hibi.yusuke@nims.go.jp and naito.masanobu@nims.go.jp

This PDF file includes:

Materials and Methods

Computational Methods

Supplementary Text

Supplementary Figures S1 to S10

Supplementary Tables S1 to S2

Captions for Data S1 to S5

References cited in this supporting information document.

Other Supplementary Materials for this manuscript include the following:

Data S1 to S5 available in [DOI: 10.26434/chemrxiv-2022-mw76d-v2](https://doi.org/10.26434/chemrxiv-2022-mw76d-v2).

Materials and Methods

Materials

Monomers were purchased from Tokyo Chemical Industry and passed through a thin pad of alumina right before polymerization. Poly(styrene-alt-methyl methacrylate) was purchased from Polymer Source (P1604-SMMAalt). For benchmark compositional analysis, living anion polymerized poly(methyl methacrylate) (M) and polystyrene (S) were purchased from Shodex ($M_n=178,00$) and Tosoh ($M_n=102,00$), respectively. Free-radical polymerized poly(ethyl methacrylate) (E) was purchased from Polymer Source. Reversible addition and fragmentation chain transfer (RAFT) agent, 2-(dodecylthiocarbonothioylthio)-2-methylpropionic acid (DDMAT), was purchased from Aldrich. Azobis(isobutyronitrile) (AIBN) was purchased from Fujifilm Wako Pure Chemical Corporation, purified by recrystallization from methanol, and stored at 0 °C. Monomers were purchased from Tokyo Chemical Industry and passed through alumina column right before use. All the other high purity chemicals were purchased from Fujifilm Wako Pure Chemical Corporation and used without further purification.

A synthetic procedure of random copolymers via free-radical copolymerization

A typical polymerization procedure is given here. In an open vial, *n*-butyl acrylate (B) (1.15 mL; 8 mmol), styrene (S) (0.23 mL; 2 mmol) and dimethyl 2,2'-azobis(isobutyrate) (23 mg; 0.1 mmol in toluene 0.1 mL) were placed. After sealing the vial with a septum,

nitrogen gas was bubbled just for one minute. The reaction mixture was stirred for an hour at 60 °C, and then dropped into methanol/water = 9/1 (v/v). The precipitated polymer was dried overnight under vacuum and subjected to pyrolysis-MS measurements for dataset preparation. High-temperature polymerization (150 °C) was conducted without addition of the initiator/toluene solution. No unexpected or unusually high safety hazards were encountered.

A synthetic procedure of B/S random copolymers via living radical copolymerization¹

To a round flask reactor, DDMAT 72.9 mg (0.2 mmol) and AIBN 9.9 mg (0.06 mmol) were placed. Styrene 2.29 mL (20 mmol), *n*-butyl acrylate 2.87 mL (20 mmol) and 1,4-dioxane (2 mL) were all combined, bubbled with nitrogen gas for 30 mins and then transfer to the degassed reactor at room temperature. The system was heated up to 70 °C to start polymerization. Small portions of the polymerization solution were taken out at 5, 10, 25, 58 h. The withdrawn solutions were subjected to ¹H NMR for conversion calculation and pyrolysis-MS for direct sequencing. The left solutions were dropped into methanol/water =9/1 (v/v). The precipitated polymers were dried under vacuum at 100 °C for 24 h and subjected to ¹H and ¹³C NMR for triad sequencing.

Sample preparation procedures for the benchmark test of the E/M/S ternary polymer film

A typical procedure for preparing ternary films is here given. A polymer mixture solution was prepared by dissolving 2.0 mg of each E/M/S in 294 mg of 1,4-dioxane (totally 2 wt% solution). The solution was then dropped in a copper pan. The solution weight was quickly measured, typically around 10 mg. The polymer weight is calculated by the solution

weight and polymer concentration, typically around 0.20 mg. The pan was left in atmosphere overnight, yielding a drop casted film for pyrolysis-MS measurement.

Pyrolysis-MS measurements

Polymer samples were pyrolyzed on a heater (ionRocket; Biochromato) at heating rate of 50 °C/min from 50 °C to 500 °C after two minutes preheating period at 50 °C. One sample thus took 12 mins measurement time. The pyrolyzed gases were continuously ionized by excited He gas with DART-ion source (DART-OS; IonSense). MS spectra were recorded with LCMS-2020 (Shimadzu) in positive-ion mode at 50 scan/min, yielding 550 spectra per a sample. The mass range was 50-1500 m/z and the interval scale was 0.05 m/z with mass resolution of 2000. The spectra were output in CDF file format, converted to Numpy format with a Python module, netCDF4. The spectra were further formatted for data size reduction as described in “Spectra formatting” section. The formatted spectral datasets used in this study are presented as Data S1-S5. All the data processing were conducted on Python3.7 on a Windows 11 laptop computer with a processing unit of AMD Ryzen9 4900HS without external GPU assistance. The total processing time depended on the dataset size and further specified mass range, which was not beyond 3 hours in this study.

Computational Methods

Mathematical notations

The notation mostly follows the common convention in signal processing. $\mathbb{R}_+^{N \times M}$ and $\mathbb{R}^{N \times M}$ represent a non-negative real matrix and real matrix, respectively, with $N \times M$ dimensions. For a matrix $\mathbf{X} \in \mathbb{R}^{N \times M}$, $\mathbf{X}_{n:} \in \mathbb{R}^{1 \times M}$, $\mathbf{X}_{:m} \in \mathbb{R}^{N \times 1}$, $X_{nm} \in \mathbb{R}$ respectively represent the n^{th} row vector, m^{th} column vector, and (n, m) -element of the matrix \mathbf{X} . \mathbf{X}^T represents a transposed matrix of \mathbf{X} . $\|\mathbf{X}\|_F$ represents Frobenius norm of \mathbf{X} . $\|\mathbf{X}_{n:}\|_1$ and $\|\mathbf{X}_{n:}\|_2$ respectively represent ℓ_1 - and ℓ_2 -norm of n^{th} row vector of \mathbf{X} . For a square matrix of \mathbf{X} , $Tr(\mathbf{X})$ represents the trace of \mathbf{X} . $\mathbf{1}_N$, $\mathbf{1}\mathbf{1}_N$ and \mathbf{I}_N respectively represent a N -dimensional all-ones vector, (N, N) -dimensional all-ones matrix and N -dimensional identity matrix.

Spectra formatting

For simplicity, in the main text, we described as if pyrolysis-MS was 1D-spectrum and dataset was $\mathbf{X} \in \mathbb{R}_+^{N \times D}$ (N : sample number, D : channel number). However, as described in pyrolysis-MS measurement section, a single pyrolysis-MS was 2D-spectrum consisting of 550 1D-spectra recorded at different temperatures. The first 100 spectra corresponded to the preheating period. The mass range was 50-1500 m/z with intervals of 0.05 m/z. The original matrix size for a single sample was thus 550×29001 . The absolute peak-intensities of DART pyrolysis-MS are sensitively affected by sample weights and ambient environments, especially humidity. As we collected data over a year in a room without controlling humidity in Japan, there should be a huge fluctuation in the relative

humidity (20-100%), making direct comparison of spectral intensities among different samples invalid. The spectral intensities were thus scaled by an internal standard peak and sample weight. As an internal standard, we used a background peak at 391 m/z contained in the first 100 spectra from preheating period, attributable to a proton adduct of di(2-ethylhexyl)phthalate, which is a plasticizer continuously released from plastics and equilibrated in a room. After scaling the spectral intensities by the internal standard peaks and sample weights, the first 100 spectra were removed from the data matrix. The spectral dimension per a sample became 450×29001 . As such high temperature resolution (1°C) is unnecessary, the spectra were averaged into T -temperature bands (typically $T \in [10, 20]$). Since the mass resolution in our setup was not high enough to precisely measure mass with 0.05 m/z accuracy, the m/z intervals were doubled to 0.1 m/z and then gaussian filter was applied; smoothing sequential 4 channels (0.4 m/z) with standard deviation of 2 channels (0.2 m/z). For further reducing data size, we specified the mass range depending on the dataset (see Table S2). In this study, we used 5 datasets: S/B binary triad sequencing (Fig. 1A), benchmark compositional analysis of E/M/S ternary polymer films (Fig. 2A), M/S binary triad sequencing (Fig. 3), M/S/B ternary triad sequencing (Fig. S8), and S/B binary pentad sequencing (Fig. 4A). For the benchmark dataset of 24 samples, the mass range was limited to 50-1000 m/z, since a region over 1000 m/z showed very subtle peaks only, deriving a 20×9501 matrix for a single sample and 460×9501 matrix for the entire dataset. For triad sequencing of S/B binary, M/S binary and M/S/B ternary, the mass range was limited to 50-500 m/z, as the fragment distribution longer than triads cannot be expressed with 5 bases (see Fig. S1). For pentad sequencing of S/B, mass range was limited to 100-670 m/z, for

covering up to pentads (the largest pentad is BBBBB at 641 m/z for proton adduct and 658 m/z for ammonium ion adduct). For sequencing, temperature range was further specified within 200-450 °C and averaged into $T=10$ (for triad sequencing) and $T=15$ (for pentad sequencing). Each spectrum was ℓ_1 -normalized and the resulting spectral dataset, $\tilde{\mathbf{X}} \in \mathbb{R}_+^{NT \times D}$ (N : sample number, T : temperature-band number, D : channel number), was subjected to the first NMF (see Fig. S4). The ℓ_1 -norms of all the NT -spectra were preserved in a diagonal matrix of $\mathbf{L}_{\tilde{\mathbf{X}}} \in \mathbb{R}_+^{NT \times NT}$ for recovering the spectral intensities after the first NMF: if the normalized $\tilde{\mathbf{X}}$ was factorized into $\tilde{\mathbf{X}} \approx \tilde{\mathbf{A}}\mathbf{S}$, then the output $\tilde{\mathbf{A}}$ was replaced with $\mathbf{L}_{\tilde{\mathbf{X}}}\tilde{\mathbf{A}}$ for recovering the absolute intensities.

Designing outline of RQMS algorithm

The RQMS algorithm is outlined in Fig. S4 as a flowchart. The RQMS algorithm is composed of three main parts: the two NMFs described in the main text and a filter for automatically identifying and removing contaminants/backgrounds fragments from the system based on canonical correlation analysis² (CCA). In the following sections, the mathematical derivations and pseudo-codes of the first NMF, CCA-filter and second NMF are described in this order.

Derivations of the first NMF

For simplicity, we wrote $\mathbf{X} \approx \mathbf{A}\mathbf{S}$ in the main text and above section, where $\mathbf{X} \in \mathbb{R}_+^{N \times D}$ and $\mathbf{A} \in \mathbb{R}_+^{N \times M}$ are sample-wise spectra and fragment abundances (FAs). However, in real implementation, $\tilde{\mathbf{X}} \in \mathbb{R}_+^{NT \times D}$ should be subjected to the first NMF rather

than $\mathbf{X} \in \mathbb{R}_+^{N \times D}$ (see Fig. S5). The first NMF is thus $\tilde{\mathbf{X}} \approx \tilde{\mathbf{A}}\mathbf{S}$ factorizing a dataset matrix $\tilde{\mathbf{X}}$ into spectrum-wise FA $\tilde{\mathbf{A}} \in \mathbb{R}_+^{NT \times M}$ and M -fragments basis spectra $\mathbf{S} \in \mathbb{R}_+^{M \times D}$. The derivation mostly followed the previously proposed NMF by Shiga *et al.*³ with ARD⁴ and soft orthogonal constraint (ARD-SO-NMF). For a while, identical and independent distributed (i.i.d.) Gaussian noise with σ^2 -variance was assumed in a probabilistic model generating $\tilde{\mathbf{X}}$ from $\tilde{\mathbf{A}}\mathbf{S}$:

$$p(\tilde{\mathbf{X}}|\tilde{\mathbf{A}}, \mathbf{S}, \sigma^2) = \prod_{i=1}^{NT} \prod_{m=1}^M \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{(\tilde{\mathbf{X}}_{im} - (\tilde{\mathbf{A}}\mathbf{S})_{im})^2}{2\sigma^2}\right\}.$$

To automatically determine the fragment number M , automatic relevance determination (ARD) based on sparse modeling was introduced, assuming exponential distribution for prior distribution of $\tilde{\mathbf{A}}$ column-wisely parameterized by $\boldsymbol{\lambda} \in \mathbb{R}^M$, i.e.:

$$p(\tilde{A}_{im}|\lambda_m) = \frac{1}{\lambda_m} \exp\left(-\frac{\tilde{A}_{im}}{\lambda_m}\right) \text{ s. t. } \lambda_m > 0, (i = 1, \dots, NT, m = 1, \dots, M),$$

$$p(\tilde{\mathbf{A}}|\boldsymbol{\lambda}) = \prod_{i=1}^{NT} \prod_{m=1}^M p(\tilde{A}_{im}|\lambda_m).$$

Note that half-gaussian distribution is also available as the prior distribution of $\tilde{\mathbf{A}}$ instead of exponential distribution, as presented in³. Assuming sparseness on $\tilde{\mathbf{A}}$ is reasonable as a certain fragment is only generated in specific samples and temperature bands. As greater λ_m tolerates greater $\|\tilde{\mathbf{A}}_{:,m}\|_1 = \sum_{i=1}^{NT} \tilde{A}_{im}$, λ_m represents the importance of the m^{th} fragment throughout the dataset. ARD starts with a large integer M and then deletes components when their λ_m becomes close to zero while updating, outputting a model with suitable component number. The total probabilistic model becomes:

$$p(\tilde{\mathbf{X}}, \tilde{\mathbf{A}}, \mathbf{S}) = p(\tilde{\mathbf{X}}|\tilde{\mathbf{A}}, \mathbf{S}, \sigma^2)p(\tilde{\mathbf{A}}|\boldsymbol{\lambda})p(\mathbf{S})p(\boldsymbol{\lambda}|a, b)$$

where uniform distribution for $p(\mathbf{S})$ on the hypersphere of $\|\mathbf{S}_m\|_2 = 1$ for $m = 1, \dots, M$, and inverse-gamma distribution for $p(\boldsymbol{\lambda}|a, b)$ parameterized by (a, b) are assumed, i.e.,

$$p(\boldsymbol{\lambda}|a, b) = \prod_{m=1}^M p(\lambda_m|a, b) = \frac{b^a}{\Gamma(a)} \lambda_m^{-(a+1)} \exp\left(-\frac{b}{\lambda_m}\right) \text{ for } m = 1, \dots, M$$

where a is a hyperparameter very close to 1 adjusting sparsity, here fixed $a = 1 + 10^{-16}$ and $\Gamma(\cdot)$ represents the gamma function. Empirically, b was determined from the relationship to expectation of \tilde{A}_{im} (13):

$$E(\tilde{A}_{im}) = \frac{b}{a - 1}. \quad (1)$$

This was further correlated to $E(X_{id})$ by

$$E(\tilde{X}_{id}) = \sum_{m=1}^M E(\tilde{A}_{im})E(S_{md}) = M \cdot E(\tilde{A}_{im})E(S_{md}).$$

By approximating $E(\tilde{X}_{id})$ as the mean of \tilde{X} (μ_X) and substituting $E(S_{md}) = \sqrt{D}$, this becomes:

$$\mu_X = M \frac{b}{(a - 1)\sqrt{D}}$$

Then, b can be determined as:

$$b = \frac{\mu_X(a - 1)\sqrt{D}}{M}. \quad (2)$$

A negative-log likelihood function now becomes

$$\begin{aligned}
L(\tilde{\mathbf{X}}, \tilde{\mathbf{A}}, \mathbf{S}, \boldsymbol{\lambda}) &= -\log[p(\tilde{\mathbf{X}}|\tilde{\mathbf{A}}, \mathbf{S})p(\tilde{\mathbf{A}}|\boldsymbol{\lambda})p(\boldsymbol{\lambda}|a, b)] \\
&= \frac{1}{2\sigma^2} \|\tilde{\mathbf{X}} - \tilde{\mathbf{A}}\mathbf{S}\|_F^2 + \frac{DNT}{2} \log 2\pi\sigma^2 + (NT + a + 1) \sum_{m=1}^M \log \lambda_m \\
&\quad + \sum_{m=1}^M \frac{1}{\lambda_m} \left(b + \sum_{i=1}^{NT} \tilde{A}_{im} \right) + M(\log \Gamma(a) - a \log b).
\end{aligned}$$

As this function is convex over $\boldsymbol{\lambda}$, the updating rule for $\boldsymbol{\lambda}$ can be determined so that $\frac{\partial L}{\partial \boldsymbol{\lambda}} \equiv \mathbf{0}$, i.e.:

$$\lambda_m = \frac{b + \sum_{i=1}^{NT} \tilde{A}_{im}}{NT + a + 1}, \text{ for } m = 1, \dots, M \quad (3)$$

The derivation so far has completely followed ARD-SO-NMF derivation reported by Shiga *et. al.*³. The i.i.d. assumption is, however, significantly violated in MS data as noise variance of a channel correlates with its signal intensity. The noise distribution does not follow Poisson distribution neither⁵. We propose to estimate a noise matrix $\mathbf{E} \in \mathbb{R}^{NT \times D}$ from the residuals of ordinary least square regression based on natural isotope peaks. The estimated noise of d^{th} channel throughout the dataset becomes:

$$\mathbf{E}_{:d} = \tilde{\mathbf{X}}_{:d} - \mathbf{M}^{(d)} (\mathbf{M}^{(d)T} \mathbf{M}^{(d)})^{-1} \mathbf{M}^{(d)T} \tilde{\mathbf{X}}_{:d}, \text{ for } d = 1, \dots, D, \quad (4)$$

where $\mathbf{M}^{(d)} = \tilde{\mathbf{X}}_{:[d-30, d-20, d-10, d+10, d+20, d+30]}$. As m/z intervals are 0.1 m/z, the channel set, $[d - 30, d - 20, d - 10, d + 10, d + 20, d + 30]$, should include ± 3 m/z isotope peaks. To intuitively understand this concept, consider a noiseless MS dataset where intensity ratio of isotopic peaks is precisely identical to natural isotope-abundance, and thus constant throughout all the spectra. The regression residuals thus become zeros. Conversely, in a real dataset, the intensity ratio of isotopic peaks is no longer constant among spectra owing to

signal noise, yielding non-zero \mathbf{E} reflecting the deviations from natural isotope-abundance ratios. A channel-wise covariance matrix $\mathbf{R} \in \mathbb{R}_+^{D \times D}$ is thus obtained by $\mathbf{R} = \frac{1}{NT} \mathbf{E}^T \mathbf{E}$. The likelihood function is then rewritten as

$$p(\tilde{\mathbf{X}}|\tilde{\mathbf{A}}, \mathbf{S}, \mathbf{R}) = \frac{1}{\sqrt{2\pi}^{-DNT} \sqrt{|\mathbf{R}|}^{NT}} \exp \left\{ -Tr \left[(\tilde{\mathbf{X}} - \tilde{\mathbf{A}}\mathbf{S})\mathbf{R}^{-1}(\tilde{\mathbf{X}} - \tilde{\mathbf{A}}\mathbf{S})^T \right] \right\}.$$

Note that \mathbf{R} is a constant matrix, not necessitating updates. The entire negative log-likelihood function now becomes

$$\begin{aligned} L(\tilde{\mathbf{X}}, \tilde{\mathbf{A}}, \mathbf{S}, \boldsymbol{\lambda}) &= -\log[p(\tilde{\mathbf{X}}|\tilde{\mathbf{A}}, \mathbf{S}, \mathbf{R})p(\tilde{\mathbf{A}}|\boldsymbol{\lambda})p(\boldsymbol{\lambda}|a, b)] \\ &= Tr \left[(\tilde{\mathbf{X}} - \tilde{\mathbf{A}}\mathbf{S})\mathbf{R}^{-1}(\tilde{\mathbf{X}} - \tilde{\mathbf{A}}\mathbf{S})^T \right] + \frac{DNT}{2} \log 2\pi + \frac{NT}{2} \log |\mathbf{R}| \\ &\quad + (NT + a + 1) \sum_{m=1}^M \log \lambda_m + \sum_{m=1}^M \frac{1}{\lambda_m} \left(b + \sum_{i=1}^{NT} \tilde{A}_{im} \right) \\ &\quad + M(\log \Gamma(a) - a \log b). \end{aligned}$$

By simplifying about $\boldsymbol{\lambda}$ using Eq. 3 and dropping constant terms, this can be rewritten as:

$$L(\tilde{\mathbf{X}}, \tilde{\mathbf{A}}, \mathbf{S}, \boldsymbol{\lambda}) = Tr \left[(\tilde{\mathbf{X}} - \tilde{\mathbf{A}}\mathbf{S})\mathbf{R}^{-1}(\tilde{\mathbf{X}} - \tilde{\mathbf{A}}\mathbf{S})^T \right] + (NT + a + 1) \sum_{m=1}^M \log \lambda_m.$$

The objective function is minimized in a framework of hierarchical alternating least square (HALS) ⁶, a vector-wise updating methods known as the fastest NMF compatible with orthogonal constraints ⁷. For convenience, the following column-vector notations are used:

$$\mathbf{a}_m \equiv \tilde{\mathbf{A}}_{:m} \in \mathbb{R}_+^{NT \times 1}, \mathbf{s}_m \equiv \mathbf{S}_{m:}^T \in \mathbb{R}_+^{D \times 1}, \text{ for } m = 1, \dots, M.$$

The key of HALS is representing the residual, $\tilde{\mathbf{X}} - \tilde{\mathbf{A}}\mathbf{S}$, as $\tilde{\mathbf{X}}^{(m)} - \mathbf{a}_m \mathbf{s}_m^T$ ($m = 1, \dots, M$), where $\tilde{\mathbf{X}}^{(m)} = \tilde{\mathbf{X}} - \tilde{\mathbf{A}}\mathbf{S} + \mathbf{a}_m \mathbf{s}_m^T$. The negative-log likelihood now can be written

separately by M -components as:

$$\begin{aligned}
L(\tilde{\mathbf{X}}, \mathbf{a}_m, \mathbf{s}_m, \lambda_m(\mathbf{a}_m)) \\
= \text{Tr} \left[(\tilde{\mathbf{X}}^{(m)} - \mathbf{a}_m \mathbf{s}_m^T) \mathbf{R}^{-1} (\tilde{\mathbf{X}}^{(m)} - \mathbf{a}_m \mathbf{s}_m^T)^T \right] + (NT + a + 1) \log \lambda_m, \\
\text{for } m = 1, \dots, M.
\end{aligned}$$

The orthogonal constraints on rows of \mathbf{S} can be included in L as a penalty term of $w_0 \xi_m \mathbf{s}_m^T \mathbf{s}^{(m)}$ where $\mathbf{s}^{(m)} \equiv \sum_{j \neq m}^M \mathbf{s}_j$. This represents non-orthogonality between the m^{th} component and all the others; ξ_m is a Lagrange multiplier when the orthogonal constraint is best-satisfied under the strict non-negative constraints, which is further softened by hyperparameter $w_0 \in [0, 1]$. Higher w_0 imposes the stronger orthogonal constraints on the rows of \mathbf{S} . Note that even if $w_0 = 1$ is given, the rows of \mathbf{S} would not be strictly orthogonal, as ξ_m has been estimated without nonnegative-constraints⁷. The full objective function to be minimized was finally derived as:

$$\begin{aligned}
L(\tilde{\mathbf{X}}, \mathbf{a}_m, \mathbf{s}_m, \lambda_m(\mathbf{a}_m)) \\
= \text{Tr} \left[(\tilde{\mathbf{X}}^{(m)} - \mathbf{a}_m \mathbf{s}_m^T) \mathbf{R}^{-1} (\tilde{\mathbf{X}}^{(m)} - \mathbf{a}_m \mathbf{s}_m^T)^T \right] + (NT + a + 1) \log \lambda_m \\
+ w_0 \xi_m \mathbf{s}_m^T \mathbf{s}^{(m)}.
\end{aligned}$$

The gradients of L over \mathbf{a}_m and \mathbf{s}_m are, respectively,

$$\begin{aligned}
\frac{\partial L}{\partial \mathbf{a}_m} &= (\mathbf{a}_m \mathbf{s}_m^T - \tilde{\mathbf{X}}^{(m)}) \mathbf{R}^{-1} \mathbf{s}_m + \frac{1}{\lambda_m} \mathbf{1}_{NT}, \\
\frac{\partial L}{\partial \mathbf{s}_m} &= \mathbf{R}^{-1} (\mathbf{s}_m \mathbf{a}_m^T - \tilde{\mathbf{X}}^{(m)T}) \mathbf{a}_m + w_0 \xi_m \mathbf{s}^{(m)}.
\end{aligned}$$

The \mathbf{a}_m and \mathbf{s}_m minimizing L are obtained by setting the gradients zeros:

$$\mathbf{a}_m = \frac{\tilde{\mathbf{X}}^{(m)} \mathbf{R}^{-1} \mathbf{s}_m - \frac{1}{\lambda_m} \mathbf{1}_{NT}}{\mathbf{s}_m^T \mathbf{R}^{-1} \mathbf{s}_m}, \quad (5)$$

$$\mathbf{s}_m = \tilde{\mathbf{X}}^{(m)T} \mathbf{a}_m - w_0 \xi_m \mathbf{R} \mathbf{s}^{(m)}. \quad (6)$$

Note that the scaling factor for \mathbf{s}_m was dropped as \mathbf{s}_m will be normalized after each update.

The updated \mathbf{a}_m and \mathbf{s}_m are projected to non-negative orthant after each update by, e.g.,

$$\mathbf{a}_m \leftarrow \frac{\mathbf{a}_m + |\mathbf{a}_m|}{2},$$

where $|\mathbf{a}_m|$ represents a vector having absolute values of \mathbf{a}_m elements. The estimation of

ξ_m is done by multiplying $\mathbf{s}^{(m)T} \mathbf{R}$ to Eq. 6 from the left side with using the strict orthogonal

constraint of $\mathbf{s}^{(m)T} \mathbf{s}_m = 0$ and $w_0 = 1$, deriving

$$-\mathbf{s}^{(m)T} \tilde{\mathbf{X}}^{(m)T} \mathbf{a}_m + \xi_m \mathbf{s}^{(m)T} \mathbf{R} \mathbf{s}^{(m)} = 0,$$

$$\xi_m = \frac{\mathbf{s}^{(m)T} \tilde{\mathbf{X}}^{(m)T} \mathbf{a}_m}{\mathbf{s}^{(m)T} \mathbf{R} \mathbf{s}^{(m)}} \quad (7)$$

Using based on Eq. 1-7, we propose the following code.

Algorithm 1: Pseudo-code of the first NMF

Input: ℓ_1 -normalized data matrix: $\tilde{\mathbf{X}} \in \mathbb{R}_+^{NT \times D}$, orthogonal constraint weight: w_0 , initial component number: M , iteration number: itr , margining threshold: $t \in [0.9, 0.99]$.

Output: spectra-wise fragment abundance: $\tilde{\mathbf{A}} \in \mathbb{R}_+^{NT \times M}$, M -fragment spectra: $\mathbf{S} \in \mathbb{R}_+^{M \times D}$.

Initialization

initialize $\tilde{\mathbf{A}}$ by Eq. 1

initialize \mathbf{S} by random selection of M -row vectors from $\tilde{\mathbf{X}}$ and ℓ_2 -normalization.

calculate b by Eq. 2

initialize λ by Eq. 3

calculate E by Eq. 4 and $R = \frac{1}{NT} E^T E$

Repeat until convergence criteria are satisfied:

for $m = 1, \dots, M$:

calculate $\mathbf{s}^{(m)} \leftarrow \sum_{j \neq m}^M \mathbf{s}_j$ and calculate ξ_m by Eq. 7

update \mathbf{s}_m by Eq. 6.

project \mathbf{s}_m to the non-negative orthant

ℓ_2 -normalize \mathbf{s}_m

for $m = 1, \dots, M$:

update \mathbf{a}_m by Eq. 5

project \mathbf{a}_m to the non-negative orthant

Merging very similar components

for $k = 1, \dots, M - 1$:

for $m = k + 1, \dots, M$:

if $\mathbf{s}_k^T \mathbf{s}_m > t$:

$\mathbf{a}_k \leftarrow \mathbf{a}_k + \mathbf{a}_m, \mathbf{a}_m \leftarrow \mathbf{0}$

if $\mathbf{a}_k^T \mathbf{a}_m > t \|\mathbf{a}_k\| \|\mathbf{a}_m\|$:

$\mathbf{s}_k \leftarrow \mathbf{s}_k + \frac{\|\mathbf{a}_m\|_1}{\|\mathbf{a}_k\|_1} \mathbf{s}_m, s = \|\mathbf{s}_k\|_2, \mathbf{s}_k \leftarrow \frac{\mathbf{s}_k}{s}$

$\mathbf{a}_k \leftarrow s \mathbf{a}_k, \mathbf{a}_m \leftarrow \mathbf{0}$

update λ by Eq. 3

return $\tilde{\mathbf{A}}$ and \mathbf{S}

Note that after each update of $\tilde{\mathbf{A}}$ and \mathbf{S} , the components having similar spectra or abundance distribution over the dataset are merged, reducing the component number down. By merging the components based on not only their spectra but also abundance distribution, a series of fragments peaks belonging to an identical codon would be unified in a single spectrum, deriving more interpretable outputs.

Derivation of the canonical correlation analysis filter (CCA-filter)

Backgrounds and/or contaminants are potentially included in the M -components output from the first NMF, which may distort the analytic results. We here describe our development of CCA-filter, automatically identifying and removing the contaminated components from the M -components. After removing M' -contaminants, $\tilde{\mathbf{A}} \in \mathbb{R}_+^{NT \times M}$ and $\mathbf{S} \in \mathbb{R}_+^{M \times D}$ are respectively reduced to $\tilde{\mathbf{A}} \in \mathbb{R}_+^{NT \times (M-M')}$ and $\mathbf{S} \in \mathbb{R}_+^{(M-M') \times D}$. For simplicity, in the main text and following second NMF derivation, M rather than $M - M'$ is consistently used for representing the component number.

To use CCA-filter, a background spectrum $\mathbf{X}_{BG} \in \mathbb{R}_+^{T \times D}$ needs to be included in a dataset before the first NMF is applied. If some contaminants are expected, e.g., residual solvents in cast-film samples, they can be mixed and measured in a single spectrum, which can be used as \mathbf{X}_{BG} . Conceptually, CCA-filter finds out M' -spectra from ℓ_2 -normalized $\mathbf{S} \in \mathbb{R}_+^{M \times D}$, of which peak patterns are also included in \mathbf{X}_{BG} . Assume we would like to examine m -th component if it is contaminant or fragment generated from the samples. The first step

is classifying M -spectral set \mathbf{S} into two groups: spectral sets $\mathbf{Y} \in \mathbb{R}_+^{M_{sim} \times D}$ and $\mathbf{Z} \in \mathbb{R}_+^{M_{dis} \times D}$, respectively similar and dissimilar to \mathbf{S}_m . Here, M_{sim} and M_{dis} are the numbers of similar and dissimilar spectra, respectively. The classification is based on cosine similarity. After the classification, the following inequalities should be hold:

$$\mathbf{S}_m: \mathbf{Y}_{m'}:^T \geq t_1, \text{ for } m' = 1, \dots, M_{sim},$$

$$\mathbf{S}_m: \mathbf{Z}_{m'}:^T < t_1, \text{ for } m' = 1, \dots, M_{dis},$$

where $t_1 \in [0, 1]$ is a threshold, in this study fixed at $t_1 = 0.2$. Note that \mathbf{S}_m is always classified into \mathbf{Y} , and put at the top row of \mathbf{Y} . The dissimilar spectral set \mathbf{Z} is further combined with \mathbf{X}_{BG} :

$$\mathbf{Z} \leftarrow \begin{pmatrix} \mathbf{Z} \\ \mathbf{X}_{BG} \end{pmatrix}.$$

\mathbf{Y} and \mathbf{Z} are then mean-subtracted:

$$\bar{\mathbf{Y}} = \mathbf{Y} \left(\mathbf{I}_D - \frac{1}{D} \mathbf{1}\mathbf{1}_D \right) \in \mathbb{R}^{M_{sim} \times D}, \quad (8)$$

$$\bar{\mathbf{Z}} = \mathbf{Z} \left(\mathbf{I}_D - \frac{1}{D} \mathbf{1}\mathbf{1}_D \right) \in \mathbb{R}^{(M_{dis}+T) \times D}. \quad (9)$$

CCA is then applied to the two spectral set. Generally, the purpose of CCA is to generate a pair of spectra most similar to each other from two different spectral sets by linear combination, here $\bar{\mathbf{Y}}$ and $\bar{\mathbf{Z}}$. Let two coefficient vectors be $\mathbf{u} \in \mathbb{R}^{M_{sim}}$ and $\mathbf{v} \in \mathbb{R}^{M_{dis}+T}$ for $\bar{\mathbf{Y}}$ and $\bar{\mathbf{Z}}$, respectively. The new mixed spectra are defined by $\mathbf{y} \equiv \mathbf{u}^T \bar{\mathbf{Y}} \in \mathbb{R}^{1 \times D}$ and $\mathbf{z} \equiv \mathbf{v}^T \bar{\mathbf{Z}} \in \mathbb{R}^{1 \times D}$. The similarity of these two spectra is evaluated with correlation coefficient:

$$\rho = \frac{\mathbf{u}^T \mathbf{V}_{yz} \mathbf{v}}{\sqrt{\mathbf{u}^T \mathbf{V}_{yy} \mathbf{u}} \sqrt{\mathbf{v}^T \mathbf{V}_{zz} \mathbf{v}}}$$

where $\mathbf{V}_{YY} = \overline{\mathbf{Y}\mathbf{Y}^T}/D, \mathbf{V}_{ZZ} = \overline{\mathbf{Z}\mathbf{Z}^T}/D, \mathbf{V}_{YZ} = \overline{\mathbf{Y}\mathbf{Z}^T}/D$. The CCA problem now becomes:

$$(\mathbf{u}^*, \mathbf{v}^*) = \underset{\mathbf{u}, \mathbf{v}}{\arg \max} \rho.$$

The optimum $(\mathbf{u}^*) \in \mathbb{R}^{M_{sim}+M_{dis}+T}$ are given as the solution of generalized eigenvalue problem:

$$\begin{pmatrix} \mathbf{0} & \mathbf{V}_{yz} \\ \mathbf{V}_{yz}^T & \mathbf{0} \end{pmatrix} \begin{pmatrix} \mathbf{U} \\ \mathbf{V} \end{pmatrix} = \begin{pmatrix} \mathbf{V}_{yy} & \mathbf{0} \\ \mathbf{0} & \mathbf{V}_{zz} \end{pmatrix} \begin{pmatrix} \mathbf{U} \\ \mathbf{V} \end{pmatrix} \begin{pmatrix} \rho_1 & & \\ & \ddots & \\ & & \rho_{M_{sim}+M_{dis}+T} \end{pmatrix} \quad (10)$$

where $\begin{pmatrix} \mathbf{U} \\ \mathbf{V} \end{pmatrix}$ consists of all the eigen column-vectors $\begin{pmatrix} \mathbf{u}^* \\ \mathbf{v}^* \end{pmatrix}$ in a descending order of the corresponding eigenvectors, i.e., $\rho_1 \geq \rho_2 \geq \dots \geq \rho_{M_{sim}+M_{dis}+T}$. Each eigen value represents the correlation coefficient of the paired spectra of \mathbf{y} and \mathbf{z} , synthesized via corresponding eigen vectors $\begin{pmatrix} \mathbf{u}^* \\ \mathbf{v}^* \end{pmatrix}$. Again, all the basis spectra similar to the examined \mathbf{S}_m are contained in \mathbf{Y} , while the background T -spectra are contained in \mathbf{Z} . If \mathbf{Y} and \mathbf{Z} can generate very similar spectra, \mathbf{S}_m is suspected of being derived from backgrounds or contaminations. Practically, we collect all the eigenvectors $\begin{pmatrix} \mathbf{u}^* \\ \mathbf{v}^* \end{pmatrix}$ of which eigenvalues satisfy $\rho > t_2$ for given $t_2 \in [0.9, 0.99]$. If the first element of \mathbf{u}^* corresponding to the coefficient of \mathbf{S}_m significantly contributes to \mathbf{u}^* , i.e., if $\frac{|u_1^*|}{\|\mathbf{u}^*\|_1} \geq t_3$ for $t_3 \in [0, 1]$, then the m -th component is judged as a background/contamination component and removed from the system. The whole process is summarized in Algorithm 2.

Algorithm 2: CCA-filter

Input: ℓ_1 -normalized basis-spectra of the 1stNMF output; $\mathbf{S} \in \mathbb{R}_+^{M \times D}$, a background spectrum; $\mathbf{X}_{BG} \in \mathbb{R}_+^{T \times D}$, thresholds: $t_1 \in [0, 1], t_2 \in [0, 1], t_3 \in [0, 1]$, in this study $(t_1, t_2, t_3) = (0.2, 0.9, 0.5)$ is consistently used.

Output: a list of the components judged as background components.

for $m = 1, \dots, M$:

for $m' = 1, \dots, M$:

 classify $\mathbf{S}_{m'}$: into \mathbf{Y} if $\mathbf{S}_m \cdot \mathbf{S}_{m'}^T \geq t_1$ else into \mathbf{Z}

$\mathbf{Z} \leftarrow \begin{pmatrix} \mathbf{Z} \\ \mathbf{X}_{BG} \end{pmatrix}$

calculate $\bar{\mathbf{Y}}$ and $\bar{\mathbf{Z}}$ by Eq. 8-9

calculate $\mathbf{V}_{YY} = \bar{\mathbf{Y}}\bar{\mathbf{Y}}^T/D, \mathbf{V}_{ZZ} = \bar{\mathbf{Z}}\bar{\mathbf{Z}}^T/D, \mathbf{V}_{YZ} = \bar{\mathbf{Y}}\bar{\mathbf{Z}}^T/D$

obtain $\mathbf{U}^* = (\mathbf{u}^*_1, \dots, \mathbf{u}^*_{M_{sim}+M_{dis}+T})$ and $(\rho_1, \dots, \rho_{M_{sim}+M_{dis}+T})$ by solving Eq. 10

find Q such that $\rho_Q \geq t_2$ and $\rho_{Q+1} < t_2$

for $q = 1, \dots, Q$:

if $\frac{|u^*_{1q}|}{\|\mathbf{u}^*_{:q}\|_1} \geq t_3$, add the component m in a background list

return background list

All the components attributed to backgrounds/contaminants are then removed from the columns of $\tilde{\mathbf{A}}$ and rows of \mathbf{S} .

Non-negative least square fitting

In the following section, we frequently use non-negative least square (NNLS) fitting. This

finds the best non-negative coefficient vector $\mathbf{x} \in \mathbb{R}_+^n$ for approximating $\mathbf{y} \in \mathbb{R}^m$ as a linear combination of n -column vectors of a constant matrix $\Phi \in \mathbb{R}^{m \times n}$, i.e.,

$$\mathbf{x}^* = \underset{\mathbf{x}}{\operatorname{argmin}} \|\mathbf{y} - \Phi \mathbf{x}\|^2, \text{ s. t. } \mathbf{x} \geq \mathbf{0}.$$

Various algorithms have been proposed to quickly solve this problem, including alternating direction multiplier methods (ADMM)⁸. We here use ADMM-NNLS proposed by Fu. *et. al.*⁹, and use a notation of $\mathbf{x}^* = \operatorname{NNLS}(\mathbf{y}, \Phi)$. For each of L -vectors set, $\mathbf{Y} = [\mathbf{y}_1, \dots, \mathbf{y}_L]$, the corresponding coefficient vector \mathbf{x}_l^* ($l = 1, \dots, L$) is separately calculated by:

$$\mathbf{x}_l^* = \operatorname{NNLS}(\mathbf{Y}_{:l}, \Phi) \text{ for } l = 1, \dots, L$$

This is simply written in a matrix form as:

$$\mathbf{X}^* = \operatorname{NNLS}(\mathbf{Y}, \Phi)$$

where $\mathbf{X}^* = [\mathbf{x}_1^*, \dots, \mathbf{x}_L^*]$ is the set of optimal non-negative coefficient vectors. Also, a similar problem with sum-to-one constraints:

$$\mathbf{x}^* = \underset{\mathbf{x}}{\operatorname{argmin}} \|\mathbf{y} - \Phi \mathbf{x}\|^2, \text{ s. t. } \mathbf{x} \geq \mathbf{0}, \mathbf{x}^T \mathbf{1}_n = 1.$$

can be solved by ADMM as proposed by Fu. *et. al.*⁹. This is called fully constrained least square (FCLS)¹⁰, whose solution is written as:

$$\mathbf{x}_l^* = \operatorname{FCLS}(\mathbf{Y}_{:l}, \Phi) \text{ for } l = 1, \dots, L,$$

or in a matrix form:

$$\mathbf{X}^* = \operatorname{FCLS}(\mathbf{Y}, \Phi).$$

Derivation of the second NMF

The first NMF outputs spectrum-wise fragment abundance (FA) $\tilde{\mathbf{A}} \in \mathbb{R}^{NT \times M}$, which is

converted to $\mathbf{A} \in \mathbb{R}_+^{N \times M}$ by sample-wise T -spectra integration (Fig. S4), as temperature-independency of FA is unnecessary for compositional analysis. NMF still should be applied to $\tilde{\mathbf{X}} \in \mathbb{R}_+^{NT \times D}$ rather than $\mathbf{X} \in \mathbb{R}_+^{N \times D}$, i.e., $\mathbf{X} \approx \mathbf{A}\mathbf{S}$ for straightforwardly deriving \mathbf{A} , as the temperature-dependent spectral modulation facilitates identifying fragment-spectra \mathbf{S} (Fig. S5). The second NMF is applied to \mathbf{A} , i.e., $\mathbf{A} \approx \mathbf{C}\mathbf{B}$. The factorized $\mathbf{C} \in \mathbb{R}_+^{N \times K}$ represents K -polymers weight fraction in N -samples, and $\mathbf{B} \in \mathbb{R}_+^{K \times M}$ represents FAs of pure K -polymers. To emphasize the row stochastic condition of \mathbf{C} , the rows of \mathbf{A} is ℓ_1 -normalized in advance ¹¹:

$$\mathbf{A} \leftarrow \mathbf{L}_A^{-1} \mathbf{A},$$

where $\mathbf{L}_A = \text{diag}(\|\mathbf{A}_{1:}\|_1, \dots, \|\mathbf{A}_{N:}\|_1)$, which is later used for determining spectral norm of basis FAs, \mathbf{B} . We evaluate the approximate residuals by Riemann metrics, i.e.,

$$D_G(\mathbf{A}|\mathbf{C}\mathbf{B}) = \text{Tr}[(\mathbf{A} - \mathbf{C}\mathbf{B})\mathbf{G}(\mathbf{A} - \mathbf{C}\mathbf{B})^T],$$

where $\mathbf{G} = \mathbf{S}\mathbf{S}^T \in \mathbb{R}_+^{M \times M}$ ¹². Note that rows of \mathbf{S} should be ℓ_2 -normalized in advance for calculating \mathbf{G} . A lower triangular matrix $\mathbf{L} \in \mathbb{R}^{M \times M}$ is obtained via Cholesky decomposition, $\mathbf{G} = \mathbf{L}\mathbf{L}^T$ ¹³. The metrics can be rewritten as:

$$D_G(\mathbf{A}|\mathbf{C}\mathbf{B}) = \text{Tr}[(\mathbf{A} - \mathbf{C}\mathbf{B})\mathbf{L}\mathbf{L}^T(\mathbf{A} - \mathbf{C}\mathbf{B})^T] = \text{Tr}[(\hat{\mathbf{A}} - \mathbf{C}\hat{\mathbf{B}})(\hat{\mathbf{A}} - \mathbf{C}\hat{\mathbf{B}})^T] = \|\hat{\mathbf{A}} - \mathbf{C}\hat{\mathbf{B}}\|_F^2$$

where $\hat{\mathbf{A}} = \mathbf{A}\mathbf{L} \in \mathbb{R}^{N \times M}$, $\hat{\mathbf{B}} = \mathbf{B}\mathbf{L} \in \mathbb{R}^{K \times M}$. Thus, conducting NMF, $\mathbf{A} \approx \mathbf{C}\mathbf{B}$, with residuals evaluation in Riemann metrics is equivalent to factorizing $\hat{\mathbf{A}} \approx \mathbf{C}\hat{\mathbf{B}}$ in Euclidean space. The optimization problem now becomes:

$$\min_{\mathbf{C}, \hat{\mathbf{B}}} \frac{1}{2} \|\hat{\mathbf{A}} - \mathbf{C}\hat{\mathbf{B}}\|_F^2 + \frac{\alpha}{2} \text{vol}(\hat{\mathbf{B}}) + \frac{\beta}{2} \text{nonorth}(\hat{\mathbf{B}}), \quad (11)$$

$$s. t. \hat{\mathbf{B}} = \mathbf{B}\mathbf{L}, \mathbf{B} \geq 0, \mathbf{C} \geq 0, \mathbf{C}\mathbf{1}_K = \mathbf{1}_N.$$

where $vol(\widehat{\mathbf{B}})$ is the volume term of the simplex spanned by row vectors of $\widehat{\mathbf{B}}$ ¹⁴, $nonorth(\widehat{\mathbf{B}})$ is a non-orthogonality term of row vectors of $\widehat{\mathbf{B}}$ ¹⁵, and $\alpha > 0, 1 > \beta > 0$ are balancing parameters. The volume-minimization term, $vol(\widehat{\mathbf{B}})$, is introduced by Fu *et al.*¹⁴.

$$vol(\widehat{\mathbf{B}}) = \log|\det(\widehat{\mathbf{B}}\widehat{\mathbf{B}}^T + \tau\mathbf{I}_K)| = \log|\det(\mathbf{H})|,$$

where $\mathbf{H}(\widehat{\mathbf{B}}) = \widehat{\mathbf{B}}\widehat{\mathbf{B}}^T + \tau\mathbf{I}_K$ and τ is a small regularization parameter, here fixed at $\tau = 10^{-8}$. This volume term ensures ground-truth identifiability under “sufficiently scattered” condition¹⁴. Roughly speaking in chemistry words, if every component has several high-purity samples in a dataset, over 50 wt% at least, the unique solution \mathbf{C}^* coincides with the true composition. As this criterion is often violated in practical datasets, we found that introducing $nonorth(\widehat{\mathbf{B}})$ would be helpful, particularly for biased dataset (Fig. S2, here, $\alpha = \beta = 0.1$). The non-orthogonality term is defined for penalizing non-diagonal non-zero elements¹⁵, i.e.:

$$nonorth(\widehat{\mathbf{B}}) = Tr\left(\Lambda(\widehat{\mathbf{B}}\widehat{\mathbf{B}}^T - diag(\widehat{\mathbf{B}}\widehat{\mathbf{B}}^T))\right),$$

where $\Lambda \in \mathbb{R}^{K \times K}$ is a symmetric Lagrange multiplier matrix. Importantly, the orthogonal constraints are imposed on $\widehat{\mathbf{B}}$ rather than \mathbf{B} , restricting heavy peak-overlapping among different polymers in the original spectral space, rational in MS. In contrast, the non-negative constraints are imposed on \mathbf{B} rather than $\widehat{\mathbf{B}}$, as \mathbf{L} and thus $\widehat{\mathbf{B}}$ may take negative elements. Fu *et al.*¹⁴ reported that volume-minimized factorization is very sensitive to outliers, and suggested to discount the importance of outliers by introducing a weight matrix¹⁴:

$$w_n = \frac{p}{2} \left(\|\widehat{\mathbf{A}}_{n\cdot} - \mathbf{C}_{n\cdot} \widehat{\mathbf{B}}\|^2 + \varepsilon \right)^{\frac{p-2}{2}} \quad (12)$$

$$\mathbf{W} = \text{diag}(w_1, \dots, w_N)$$

where $p \in (0, 2]$ and ε is a small regularization positive number. Smaller p more strongly discounts the samples with large fitting error. In this work, $p \in \{1, 1.5\}$ and $\varepsilon = 10^{-8}$ were used. With \mathbf{W} , the optimization problem is rewritten as:

$$\min_{\mathbf{C}, \widehat{\mathbf{B}}} \frac{1}{2} \text{Tr} \left[\mathbf{W} (\widehat{\mathbf{A}} - \mathbf{C} \widehat{\mathbf{B}}) (\widehat{\mathbf{A}} - \mathbf{C} \widehat{\mathbf{B}})^T \right] + \frac{\alpha}{2} \text{vol}(\widehat{\mathbf{B}}) + \frac{\beta}{2} \text{nonorth}(\widehat{\mathbf{B}}),$$

$$s. t. \widehat{\mathbf{B}} = \mathbf{B} \mathbf{L}, \mathbf{B} \geq 0, \mathbf{C} \geq 0, \mathbf{C} \mathbf{1}_K = \mathbf{1}_N.$$

We solved this optimization by updating \mathbf{C} and $\widehat{\mathbf{B}}$ alternatively based on block coordinate descent (BCD) theory¹⁶. Assume we have $(\mathbf{C}^{(t)}, \widehat{\mathbf{B}}^{(t)})$ after t^{th} updating. The initial $\widehat{\mathbf{B}}^{(0)}$ can be selected out from the N -rows of $\widehat{\mathbf{A}}$ via vertex component analysis (VCA)¹⁷. As the penalty terms are not related to \mathbf{C} , updating \mathbf{C} based on a fixed $\widehat{\mathbf{B}}^{(t)}$ is straightforward by using the forementioned ADMM-FCLS, i.e.:

$$\mathbf{C}^{(t+1)T} = \text{FCLS} \left(\widehat{\mathbf{A}}^T, \widehat{\mathbf{B}}^{(t)T} \right). \quad (13)$$

Note that \mathbf{W} is a diagonal matrix and thus negligible for updating \mathbf{C} . Now we focus on updating $\widehat{\mathbf{B}}$ based on the fixed $\mathbf{C}^{(t)}$. To efficiently obtain the updating rule from $\widehat{\mathbf{B}}^{(t)}$ to $\widehat{\mathbf{B}}^{(t+1)}$, majorizer function of $\text{vol}(\widehat{\mathbf{B}})$ is introduced. After t^{th} update, based on fixed $\mathbf{H}^{(t)} = \mathbf{H}(\widehat{\mathbf{B}}^{(t)})$, the following tangent inequality holds:

$$\log|\det(\mathbf{H})| \leq \log|\det(\mathbf{H}^{(t)})| + \text{Tr} \left[(\nabla_{\mathbf{H}^{(t)}} \log|\det(\mathbf{H})|)^T (\mathbf{H} - \mathbf{H}^{(t)}) \right],$$

here $\nabla_{\mathbf{H}^{(t)}} \log|\det(\mathbf{H})|$ is the gradient of $\log|\det(\mathbf{H})|$ over \mathbf{H} at $\mathbf{H}^{(t)}$. As

$\nabla_{\mathbf{H}^{(t)}} \log|\det(\mathbf{H})| = \mathbf{H}^{(t)-T}$, $vol(\widehat{\mathbf{B}})$ can be majorized as:

$$vol(\widehat{\mathbf{B}}) = \log|\det(\mathbf{H})| \leq Tr[\mathbf{H}^{(t)-1}\mathbf{H}] + const = Tr[\mathbf{F}^{(t)}\widehat{\mathbf{B}}\widehat{\mathbf{B}}^T] + const,$$

where $\mathbf{F}^{(t)} = \mathbf{H}^{(t)-1}$ and $const$ is a constant term not related to $\widehat{\mathbf{B}}$. After replacing $vol(\widehat{\mathbf{B}})$ with the majorizer $Tr[\mathbf{F}\widehat{\mathbf{B}}\widehat{\mathbf{B}}^T]$, all the penalty terms are combined as:

$$\begin{aligned} \frac{\alpha}{2} vol(\widehat{\mathbf{B}}) + \frac{\beta}{2} nonorth(\widehat{\mathbf{B}}) &\leq \frac{\alpha}{2} Tr(\mathbf{F}^{(t)}\widehat{\mathbf{B}}\widehat{\mathbf{B}}^T) + \frac{\beta}{2} Tr(\mathbf{\Lambda}^{(t)}\widehat{\mathbf{B}}\widehat{\mathbf{B}}^T) + const \\ &= \frac{1}{2} Tr(\mathbf{V}\widehat{\mathbf{B}}\widehat{\mathbf{B}}^T) + const, \end{aligned}$$

where $\mathbf{V} = \alpha\mathbf{F}^{(t)} + \beta\mathbf{\Lambda}^{(t)}$. $\widehat{\mathbf{B}}^{(t+1)}$ is then updated by solving the following problem:

$$\begin{aligned} \widehat{\mathbf{B}}^{(t+1)} &= \arg \min_{\widehat{\mathbf{B}}} \frac{1}{2} Tr[\mathbf{W}(\widehat{\mathbf{A}} - \mathbf{C}^{(t)}\widehat{\mathbf{B}})(\widehat{\mathbf{A}} - \mathbf{C}^{(t)}\widehat{\mathbf{B}})^T] + \frac{1}{2} Tr(\mathbf{V}\widehat{\mathbf{B}}\widehat{\mathbf{B}}^T) \quad (14) \\ &s. t. \widehat{\mathbf{B}} = \mathbf{B}\mathbf{L}, \mathbf{B} \geq 0. \end{aligned}$$

To simplify the complicated constraints, $\widehat{\mathbf{B}} = \mathbf{B}\mathbf{L}$ is embedded into the objective function by using Lagrangian multiplier $\mathbf{Z} \in \mathbb{R}^{K \times M}$ in a framework of ADMM:

$$\begin{aligned} (\mathbf{B}^{(t+1)}, \widehat{\mathbf{B}}^{(t+1)}) &= \arg \min_{\mathbf{B}, \widehat{\mathbf{B}}} \max_{\mathbf{Z}} \left\{ f(\widehat{\mathbf{B}}) + Tr[\mathbf{Z}^T(\widehat{\mathbf{B}} - \mathbf{B}\mathbf{L})] - \frac{1}{2\mu} \|\mathbf{Z} - \mathbf{Z}'\|_F^2 \right\}, \\ &s. t. \mathbf{B} \geq 0, \end{aligned}$$

where $f(\widehat{\mathbf{B}}) \equiv \frac{1}{2} Tr[\mathbf{W}(\widehat{\mathbf{A}} - \mathbf{C}^{(t)}\widehat{\mathbf{B}})(\widehat{\mathbf{A}} - \mathbf{C}^{(t)}\widehat{\mathbf{B}})^T] + \frac{1}{2} Tr(\mathbf{V}\widehat{\mathbf{B}}\widehat{\mathbf{B}}^T)$ and μ is a hyperparameter for ADMM (in this work, $\mu = 1$ is consistently used.). Here, \mathbf{Z}' represents \mathbf{Z} in the previous cycle. The objective function is maximized over \mathbf{Z} when

$$\mathbf{Z} = \mathbf{Z}' + \mu(\widehat{\mathbf{B}} - \mathbf{B}\mathbf{L}). \quad (15)$$

The maximized objective function over \mathbf{Z} now becomes

$$g(\mathbf{B}, \widehat{\mathbf{B}}; \mathbf{Z}) \equiv f(\widehat{\mathbf{B}}) + \frac{\mu}{2} \left\| \widehat{\mathbf{B}} - \mathbf{B}\mathbf{L} + \frac{1}{\mu} \mathbf{Z} \right\|_F^2, \mathbf{B} \geq 0,$$

which is then minimized over $(\mathbf{B}, \widehat{\mathbf{B}})$. $(\mathbf{B}, \widehat{\mathbf{B}}, \mathbf{Z})$ is thus cyclically updated via Algorithm3.

Algorithm 2: ADMM for solving the optimization of Eq. 14

input: $\mathbf{B}^{(t)}, \widehat{\mathbf{B}}^{(t)}$, hyperparameter μ ($\mu = 1$ consistently used in this study), function $g(\mathbf{B}, \widehat{\mathbf{B}})$

output: $\mathbf{B}^{(t+1)}, \widehat{\mathbf{B}}^{(t+1)}$

initialize: $q = 0, \mathbf{B}_q \leftarrow \mathbf{B}^{(t)}, \widehat{\mathbf{B}}_q \leftarrow \widehat{\mathbf{B}}^{(t)}, \mathbf{Z}_q \leftarrow \mathbf{0}$

repeat until convergence:

$$\widehat{\mathbf{B}}_{q+1} = \underset{\widehat{\mathbf{B}}}{\operatorname{argmin}} g(\widehat{\mathbf{B}}; \mathbf{B}_q, \mathbf{Z}_q), \text{ solved by Eq. 16 as described below}$$

$$\mathbf{B}_{q+1} = \underset{\mathbf{B} \geq \mathbf{0}}{\operatorname{argmin}} g(\mathbf{B}; \widehat{\mathbf{B}}_{q+1}, \mathbf{Z}_q), \text{ solved by Eq. 17 as described below}$$

$$\mathbf{Z}_{q+1} = \mathbf{Z}_q + \mu(\widehat{\mathbf{B}}_{q+1} - \mathbf{B}_{q+1}\mathbf{L}) \text{ (Eq. 15)}$$

$$q \leftarrow q + 1$$

$$\mathbf{B}^{(t+1)} \leftarrow \mathbf{B}_q, \widehat{\mathbf{B}}^{(t+1)} \leftarrow \widehat{\mathbf{B}}_q.$$

Return $\mathbf{B}^{(t+1)}, \widehat{\mathbf{B}}^{(t+1)}$

As the objective function $g(\widehat{\mathbf{B}}; \mathbf{B}_q, \mathbf{Z}_q)$ has a simple quadratic form without constraints about $\widehat{\mathbf{B}}$, $\widehat{\mathbf{B}}_{q+1}$ can be determined in a closed form by setting the derivative of g zero, i.e.:

$$\frac{\partial g}{\partial \widehat{\mathbf{B}}} = \left(\mathbf{C}^{(t)T} \mathbf{W} \mathbf{C}^{(t+1)} + \mathbf{V} + \mu \mathbf{I}_K \right) \widehat{\mathbf{B}} - \mathbf{C}^{(t)T} \mathbf{W} \widehat{\mathbf{A}} - \mathbf{B}_q \mathbf{L} + \frac{1}{\mu} \mathbf{Z} \equiv \mathbf{0}$$

$$\widehat{\mathbf{B}}_{q+1} = \left(\mathbf{C}^{(t)T} \mathbf{W} \mathbf{C}^{(t)} + \mathbf{V} + \mu \mathbf{I}_K \right)^{-1} \left(\mathbf{C}^{(t)T} \mathbf{W} \widehat{\mathbf{A}} + \mathbf{B}_q \mathbf{L} - \frac{1}{\mu} \mathbf{Z} \right). \quad (16)$$

As $g(\mathbf{B}; \widehat{\mathbf{B}}_{q+1}, \mathbf{Z}_q)$ has a simple quadratic form about \mathbf{B} with non-negative constrains, the NNLS algorithm can be directly applied, i.e.:

$$\mathbf{B}_{q+1}^T = \text{NNLS} \left(\left(\widehat{\mathbf{B}}_{q+1} + \frac{1}{\mu} \mathbf{Z} \right)^T, \mathbf{L}^T \right). \quad (17)$$

Updating rules for $(\widehat{\mathbf{B}}, \mathbf{B}, \mathbf{C})$ to solve the original optimization problem (Eq. 11) have been obtained so far, but updating rule for the multiplier matrix $\mathbf{V} = \alpha \mathbf{F}^{(t)} + \beta \boldsymbol{\Lambda}^{(t)}$ remains unknown. $\mathbf{F}^{(t)}$ is directly calculated from $\widehat{\mathbf{B}}^{(t)}$:

$$\mathbf{F}^{(t)} = \left(\widehat{\mathbf{B}}^{(t)} \widehat{\mathbf{B}}^{(t)T} + \tau \mathbf{I}_K \right)^{-1}. \quad (18)$$

$\boldsymbol{\Lambda}^{(t)}$ is estimated from strict orthogonal conditions, i.e., $\widehat{\mathbf{B}} \widehat{\mathbf{B}}^T = \text{diag}(\widehat{\mathbf{B}} \widehat{\mathbf{B}}^T) \equiv \mathbf{D}$, $\beta = 1$ and

$$\frac{\partial f(\widehat{\mathbf{B}}; \mathbf{C}^{(t)})}{\partial \widehat{\mathbf{B}}} = \mathbf{C}^{(t)T} \mathbf{W} (\mathbf{C}^{(t)} \widehat{\mathbf{B}} - \widehat{\mathbf{A}}) + (\alpha \mathbf{F}^{(t)} + \boldsymbol{\Lambda}^{(t)}) \widehat{\mathbf{B}} \equiv 0. \quad (19)$$

The multiplier matrix $\boldsymbol{\Lambda}^{(t)}$ at the strict condition is obtained by multiplying $\widehat{\mathbf{B}}^T$ from the right side of Eq. 19:

$$\begin{aligned} \mathbf{C}^{(t)T} \mathbf{W} (\mathbf{C}^{(t)} \mathbf{D} - \widehat{\mathbf{A}} \widehat{\mathbf{B}}^T) + (\alpha \mathbf{F}^{(t)} + \boldsymbol{\Lambda}^{(t)}) \mathbf{D} &\equiv 0, \\ \boldsymbol{\Lambda}^{(t)} &= \mathbf{C}^{(t)T} \mathbf{W} (\widehat{\mathbf{A}} \widehat{\mathbf{B}}^T \mathbf{D}^{-1} - \mathbf{C}^{(t)}) - \alpha \mathbf{F}^{(t)}. \end{aligned} \quad (20)$$

By combining Eq. 18 and Eq. 20, updating rule of \mathbf{V} is obtained as:

$$\mathbf{V} = \alpha (1 - \beta) \left(\widehat{\mathbf{B}}^{(t)} \widehat{\mathbf{B}}^{(t)T} + \tau \mathbf{I}_K \right)^{-1} + \beta \mathbf{C}^{(t)T} \mathbf{W} (\widehat{\mathbf{A}} \widehat{\mathbf{B}}^T \mathbf{D}^{-1} - \mathbf{C}^{(t)}). \quad (21)$$

The volume penalty term is interpretable as a shrinking force imposed on the simplex

spanned by the rows of $\widehat{\mathbf{B}}$ ¹⁸. Conversely, the orthogonal term is an expanding force, as evident from Eq. 21; the weight $\beta \in [0, 1]$ buffers the first term of volume-minimization. As the residual-minimization term, $\|\widehat{\mathbf{A}} - \mathbf{C}\widehat{\mathbf{B}}\|_F^2$, is also interpretable as expanding force¹⁸, solving the original problem of Eq. 11 is conceptually identical to finding the equilibrium point of the three expanding and shrinking forces imposed on the simplex. Based on the updating rules, we propose the following algorithm to solve the second NMF.

Algorithm 3: Pseudo-code for the second NMF for solving Eq. 11

Input: output of the first NMF (sample-wise FA: $\mathbf{A} \in \mathbb{R}_+^{N \times M}$, fragment spectra: $\mathbf{S} \in \mathbb{R}_+^{M \times D}$), basis polymer number: K , weights for penalty terms: (α, β) , weight for outliers: p

Output: polymer weight fraction: $\mathbf{C} \in \mathbb{R}_+^{N \times K}$, FAs of basis polymers per unit weight: $\mathbf{B} \in \mathbb{R}_+^{K \times M}$

Initialization

calculate \mathbf{L} via Cholesky decomposition of $\mathbf{S}\mathbf{S}^T$

set $\widehat{\mathbf{A}} = \mathbf{A}\mathbf{L}$

initialize $\widehat{\mathbf{B}}^{(0)}$ by selecting out K -rows from $\widehat{\mathbf{A}}$ via VCA algorithm

set $\mathbf{B}^{(0)} = \widehat{\mathbf{B}}^{(0)}\mathbf{L}^{-1}$

initialize $\mathbf{C}^{(0)}$ by $\mathbf{C}^{(0)T} = FCLS(\widehat{\mathbf{A}}^T, \widehat{\mathbf{B}}^{(0)T})$

initialize \mathbf{W} by Eq. 12

initialize \mathbf{V} based on $\mathbf{B}^{(0)}$ and $\mathbf{C}^{(0)}$ by Eq. 21

Repeat until convergence criteria is satisfied:

update $\widehat{\mathbf{B}}, \mathbf{B}$ by algorithm 3

update \mathbf{C} by Eq. 13

update \mathbf{W}, \mathbf{V} by Eq. 12 and S18, respectively

return $\mathbf{C} \in \mathbb{R}_+^{N \times K}$ and \mathbf{B}

As row vectors of \mathbf{A} have been ℓ_1 -normalized in advance, the factorized \mathbf{B} are also row stochastic in noiseless cases¹¹. We still need to recover the norms of row vectors of \mathbf{B} based on \mathbf{L}_A , which are the preserved norms of \mathbf{A} rows before being normalized. For the norm-recovered $\mathbf{L}_A \mathbf{A}$, the following equation holds:

$$\mathbf{L}_A \mathbf{A} \approx \mathbf{L}_A \mathbf{C} \mathbf{B} = \mathbf{L}_A \mathbf{C} (\mathbf{L}_{B'}^{-1} \mathbf{B}') = (\mathbf{L}_{\tilde{A}} \mathbf{C} \mathbf{L}_{B'}^{-1}) \mathbf{B}' = \mathbf{C}' \mathbf{B}',$$

where $\mathbf{B}' = \mathbf{L}_{B'} \mathbf{B}$ is the norm-recovered variant of \mathbf{B} , $\mathbf{L}_{B'} = \text{diag}(\|\mathbf{B}'_{1:}\|_1, \dots, \|\mathbf{B}'_{N:}\|_1)$ is the spectral norm of the basis spectral, and $\mathbf{C}' = \mathbf{L}_{\tilde{A}} \mathbf{C} \mathbf{L}_{B'}^{-1}$ is the fraction of the norm-recovered basis \mathbf{B}' . Because the row-stochastic conditions of \mathbf{C}' still should be hold if the spectral intensities of raw spectra $\tilde{\mathbf{X}}$ have been properly corrected (see the section of “Spectra Formatting”), $\mathbf{L}_{B'}$ can be determined so that $\mathbf{C}' \mathbf{1}_K = \mathbf{1}_N$ would be best satisfied. For sequencing, monomer composition ratio would be also useful to determine $\mathbf{L}_{B'}$ particularly when the intensity correction of $\tilde{\mathbf{X}}$ is difficult and row-stochastic conditions of \mathbf{C}' are unreliable. The reconstructed basis spectra \mathbf{P} with absolute intensities in original D -dimensional spectral space thus can be obtained as $\mathbf{P} = \mathbf{B}' \mathbf{S}$ (simply written as $\mathbf{P} = \mathbf{B} \mathbf{S}$ in the main text), which showed good consistency both in spectral shapes and absolute intensities with the observed basis spectra in benchmark test (Fig 2A).

Sequential projection of a target pyrolysis-MS spectrum

Here, we describe how to project a targeted 2D MS spectrum, $\tilde{\mathbf{X}}_t \in \mathbb{R}_+^{T \times D}$ of a target copolymer onto the subspaces spanned by the rows of $\mathbf{S} \in \mathbb{R}_+^{M \times D}$ and subsequently $\mathbf{B} \in \mathbb{R}_+^{K \times M}$. The first projection onto \mathbf{S} -subspace can be simply done via NNLS fitting:

$$\tilde{\mathbf{A}}_t^T = \text{NNLS}\left(\mathbf{R}^{-\frac{1}{2}}\tilde{\mathbf{X}}_t^T, \mathbf{R}^{-\frac{1}{2}}\mathbf{S}^T\right),$$

where noise covariance $\mathbf{R} = \frac{1}{NT}\mathbf{E}^T\mathbf{E} \in \mathbb{R}^{D \times D}$ adjusts channel-wise importance based on a noise-matrix $\mathbf{E} \in \mathbb{R}^{NT \times D}$ estimated from isotope peaks. The spectrum-wise FA $\tilde{\mathbf{A}}_t$ is integrated along the temperature axis to be converted into $\mathbf{A}_t \in \mathbb{R}_+^{1 \times M}$, which is subsequently projected onto \mathbf{B} -subspace:

$$\mathbf{C}_t^T = \text{NNLS}(\mathbf{L}^T\mathbf{A}_t^T, \mathbf{L}^T\mathbf{B}^T),$$

$$\mathbf{C}_t \leftarrow \frac{\mathbf{C}_t}{\|\mathbf{C}_t\|_1},$$

where $\mathbf{C}_t \in \mathbb{R}_+^{1 \times K}$ is the fraction of K -basis polymers and $\mathbf{L} \in \mathbb{R}^{M \times M}$ is obtained by Cholesky decomposition of $\mathbf{S}\mathbf{S}^T$. Importantly, \mathbf{C}_t is scale-invariant owing to ℓ_1 -normalization. This means the weight of the target polymer subjected to pyrolysis-MS is unnecessary, allowing direct sequencing without chemical purification. In contrast, sample weights included in dataset and used for learning \mathbf{S} and \mathbf{B} should be accurately measured, as they influence the absolute intensities of the estimated basis spectra, i.e., $\mathbf{P} = \mathbf{B}\mathbf{S} \in \mathbb{R}_+^{K \times D}$.

Prediction of sequence distribution from the monomer reactivity ratio

We here use the following notations; monomer reactivity ratios for monomer 1 and monomer 2: r_1 and r_2 ; molar fraction of monomer 1 in polymerization solution: f ; feed molar fraction: f^0 ; molar fraction of monomer 1 in instantaneously generated polymers: F ; unpolymerized total monomer concentration: $[M]$; initial monomer concentration: $[M]_0$; monomer conversion: $conv = 1 - \frac{[M]}{[M]_0}$. As well known, Lewis-Mayo equation¹⁹ predicts monomer fraction in instantaneously generated polymers:

$$\frac{F}{1-F} = \frac{f}{1-f} \frac{r_1 f + (1-f)}{f + r_2(1-f)}$$

By solving the equation, F can be written as a function of f under given r_1, r_2 :

$$F = \frac{f + (r_1 - 1)f^2}{(r_1 + r_2 - 2)f^2 + 2(1 - r_2)f + r_2} \equiv F(f; r_1, r_2).$$

Skeist correlated monomer conversion and f via the following equation²⁰:

$$\log \frac{[M]}{[M]_0} = \int_{f^0}^f \frac{df}{F(f; r_1, r_2) - f}$$

As the right hand is all about f , the integral can be numerically calculated from a given f^0 to f , allowing us to calculate $conv$ given f . Next, we calculate the triad fraction based on r_1, r_2 and f . According to the first order Markov terminal model²¹, triad fraction in instantaneously generated (or elongated in living polymerization) polymers, $\{T_{111}, T_{112}, T_{212}, T_{222}, T_{221}, T_{121}\}$, is determined by the following equations:

$$T_{111} = (1 - P_{1 \rightarrow 2})^2,$$

$$T_{112} = 2P_{1\rightarrow 2}(1 - P_{1\rightarrow 2}),$$

$$T_{212} = (P_{1\rightarrow 2})^2,$$

$$T_{222} = (1 - P_{2\rightarrow 1})^2,$$

$$T_{221} = 2P_{2\rightarrow 1}(1 - P_{2\rightarrow 1}),$$

$$T_{121} = (P_{2\rightarrow 1})^2,$$

where $P_{1\rightarrow 2} = \frac{1}{1+r_1f/(1-f)}$ is the probability of monomer 2 addition to monomer 1 radical,

and $P_{2\rightarrow 1} = \frac{1}{1+r_2(1-f)/f}$ is the probability of monomer 1 addition to monomer 2 radical.

After separately calculating *conv* and $\{T_{111}, T_{112}, T_{212}, T_{222}, T_{221}, T_{121}\}$ based on given (f, r_1, r_2) , we plotted $\{T_{111}, T_{112}, T_{212}, T_{222}, T_{221}, T_{121}\}$ as a function of *conv*, as shown in Fig. S9. Note that what we can observe is accumulated triad fraction, not instantaneously generated triad fraction. We thus further integrated the instantaneously generated triad fraction along conversion, so that we can compare the predicted and observed accumulated triad fraction in Fig. 4C. To predict the triad fraction of S/B copolymers synthesized at 70 °C (Fig. 4C), we used the reported value $(r_S, r_B) = (0.70, 0.17)$ at 60 °C²².

Determination of the S/B triad fraction via ¹H NMR

As shown in Fig. S10C, -O-CH₂- protons of B-units around 4.1-3.5 ppm are sensitive to the adjacent monomer species, splitting into three major peaks, {X, Y, Z} according to the literature²¹. The triad fraction and peak fraction are correlated by the following equations:

$$\begin{pmatrix} T_{BBB} \\ T_{BBS} \\ T_{SBS} \end{pmatrix} = \begin{pmatrix} 1 & 1 - \sigma & (1 - \sigma)^2 \\ 0 & \sigma & 2\sigma(1 - \sigma) \\ 0 & 0 & \sigma^2 \end{pmatrix}^{-1} \begin{pmatrix} x \\ y \\ z \end{pmatrix},$$

where (x, y, z) is the peak-area fraction of $\{X, Y, Z\}$ and $\sigma = 0.85$ is a coisotactic parameter²¹. For determining S-centered triad fraction, ^{13}C NMR is sometimes available. The quaternary carbon of the benzene ring is known for being sequence-sensitive. However, as shown in Fig. S10B, SSS, SSB, BSB peaks significantly overlap with some hardly attributable peaks, hindering quantitative analysis. We thus used only B-centered triad fraction determined from ^1H NMR, for verifying our RQMS pentad sequencing (Fig. 4B-C). Also see the following section for converting RQMS pentad-fractions to B-centered triad fraction so that NMR and RQMS are comparable.

Downgrading the S/B pentad-fraction.

As NMR is not sensitive enough to determine S/B pentad distribution, we downgraded the RQMS pentad-sequence results to B-centered triad fractions so that RQMS and NMR results are comparable. We first prepared a matrix $\mathbf{T}_B \in \mathbb{R}_+^{9 \times 3}$ connecting nine basis sequence-defined copolymers to B-centered 3 triads abundances, as shown in Table S2. We then obtained pentad sequencing results, $\mathbf{C}_t \in \mathbb{R}_+^{1 \times K}$ (here $K = 9$) via procedure described in the section of “*Sequential projection of a target pyrolysis-MS spectrum onto the learned subspaces*”. The B-centered triad fractions can be obtained as $\mathbf{C}_t \mathbf{T}_B$ after normalization so that $\mathbf{C}_t \mathbf{T}_B$ satisfies row-stochastic conditions. S-centered triads fraction can be similarly calculated as well with \mathbf{T}_S (table S2), which is not used in this study since S-centered triads fraction is not accessible via NMR analysis as mentioned above.

Supplementary Text

Definition of sequencing

In the main text, “sequencing” is defined sequence distribution analysis. This definition could be controversial, since Lutz *et al.* suggested to distinguish “sequencing” and “sequence analysis”²³. According to their definition, “sequencing” refers to determining uniformly-defined sequence of a monodispersed polymer ensemble, while “sequence analysis” refers to analyzing sequence tendencies of a polydisperse ensemble. The classification is thus based on the nature of the objective copolymer ensembles. As our proposed method is applicable to either mono/poly-dispersed ensembles, we do not distinguish sequencing and sequence distribution analysis in our manuscript. The sequence-defined copolymers would span the probability-simplex and thus be located at the vertices, while polydisperse copolymers would be located at interior points of the simplex of which coordinates represent the codon composition, as described in the main text. Note that we here assume periodically sequenced copolymers as sequence-defined copolymers. We do not assume arbitrarily yet uniformly sequenced copolymers as an target copolymer ensemble, since such ultimate sequence-controlled copolymerization has not yet been developed except for unpractical single-monomer-addition “oligomerization” strategy, e.g., see²⁴.

Supplementary Figures

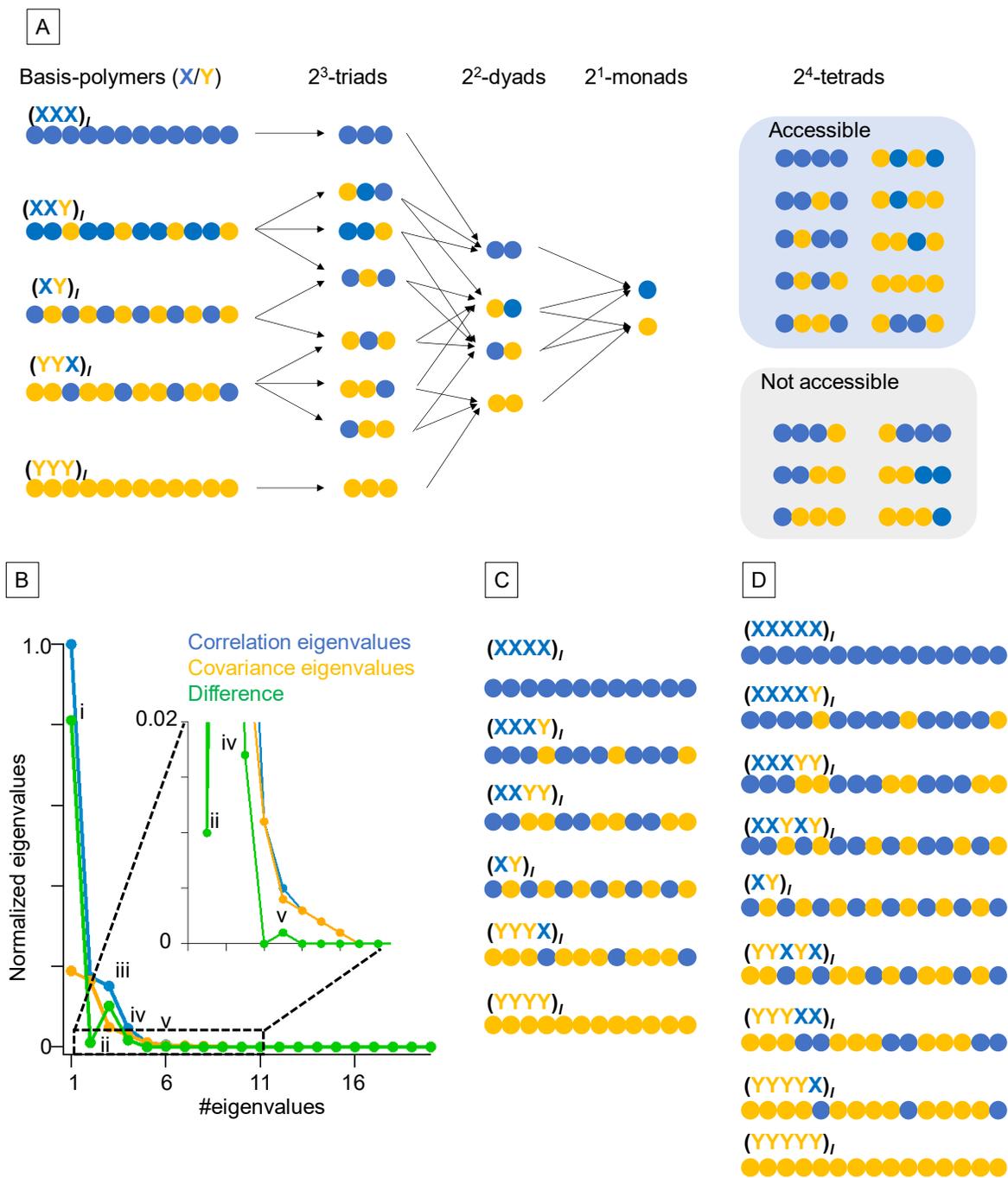


Fig. S1. Assumed basis copolymers for RQMS sequencing. (A) Five-basis copolymers for binary triad sequencing. All the 2^3 -triads are included in the five basis polymers. Pyrolysis

randomly cleaves the polymer chains generating various length of codons not limited to triad. Nevertheless, as the generation probability of a shorter codon mainly depends on the structure of the one-unit longer codon, the abundances of shorter codons than triad can be written in hierarchical chain structure. This means there are only five-basis patterns of triads-dyads-monads peak series with fixed intensity ratios, whose linear combination can express any triads-dyads-monads distributions generated from arbitrarily sequenced copolymers. To intuitively understand this, it would be helpful to consider the tetrads included in the five-basis copolymers. There are 2^4 -tetrads, six of which are not included in the five basis polymers. Therefore, when the basis number K is set five, accessible codon length is limited to triads. To express all the tetrad possibilities, we need a different basis set of six copolymers as shown in (C). In this way, we can suppress rapid increase in K for longer codon analysis though the possible number of codon sequences exponentially increases. Note that XXY and YXX are different triads which potentially generate different fragments patterns, but can be summarized in a single basis-copolymer $(XXY)_l$, because the abundance ratio of XXY and YXX is identical in any copolymers. Importantly, setting $K = 5$ is not necessarily appropriate for triad sequencing. For example, consider a monomer combination with monomer reactivity ratios of $r_X = 0$ and $r_Y > 0$. In this case, the possible triads are limited to $(XY)_l$, $(YYX)_l$, and $(YYY)_l$, as XX dyad is forbidden. K thus should be three, not five. Significantly, our developed sequence necessitates no prior knowledge nor assumption about the chemical structures of the basis copolymers. Our sequencer only necessitates the number of basis copolymers; therefore, statistical determination of the basis number K from the spectral dataset could be more effective as shown in (B). As well known, K can be determined based

on the differences between the eigenvalues of correlation and covariance matrices; the number of the non-zero difference components, marked as $i \sim v$, corresponds to K^{25} . Here, the dataset for S/B triad sequencing (Fig. 1A) was used as an example, clearly indicating that K should be set as five. **(C)** The six-basis copolymers for tetrad sequencing. **(D)** The nine-basis copolymers for tetrad sequencing.

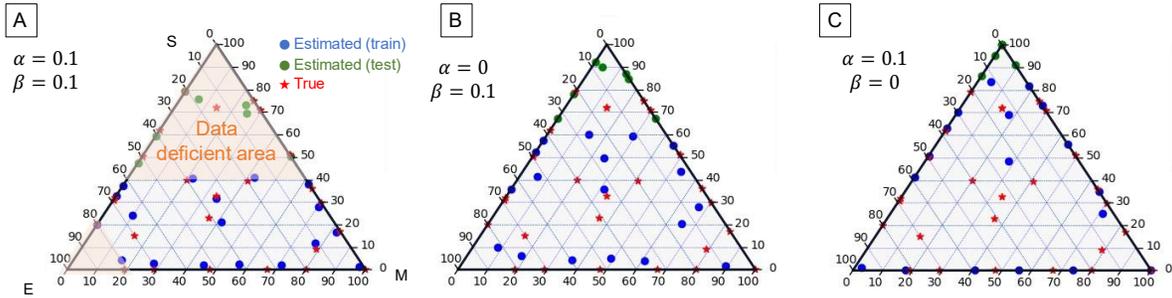


Fig. S2. Benchmark compositional analysis of the M/S/E ternary films based on a biased dataset lacking S-rich samples. The S-rich datapoints beyond 40 wt% S-fraction (green points) were not used for learning **S** and **B**, but were sequentially projected onto the subspaces spanned by **S** and **B** pre-learned from blue datapoints. The data deficient areas are represented by orange triangles. The weighting parameters, α and β , for balancing the penalty terms of volume-minimization and orthogonal constraints in the second NMF (see Eq. 11 in “Derivation of the second NMF” and the related sections) are (A) $\alpha = \beta = 0.1$, (B) $\alpha = 0, \beta = 0.1$ and (C) $\alpha = 0.1, \beta = 0$. Only when both constraints are active (A), RQMS algorithm output accurate result.

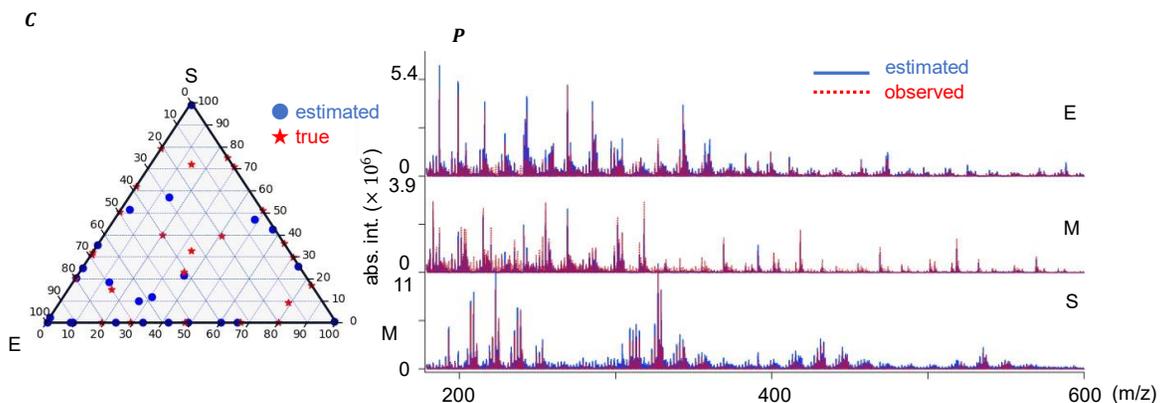


Fig. S3. Benchmark compositional analyses of E/M/S ternary films via single-step NMF, $X \approx CP$, with $K = 3, \alpha = \beta = 0.1$ (see details of these hyper parameters in the derivation of the second NMF section). The estimation errors of both C and P became much greater as compared to Fig. 2A, suggesting the effectiveness of the two-step NMF for RQMS.

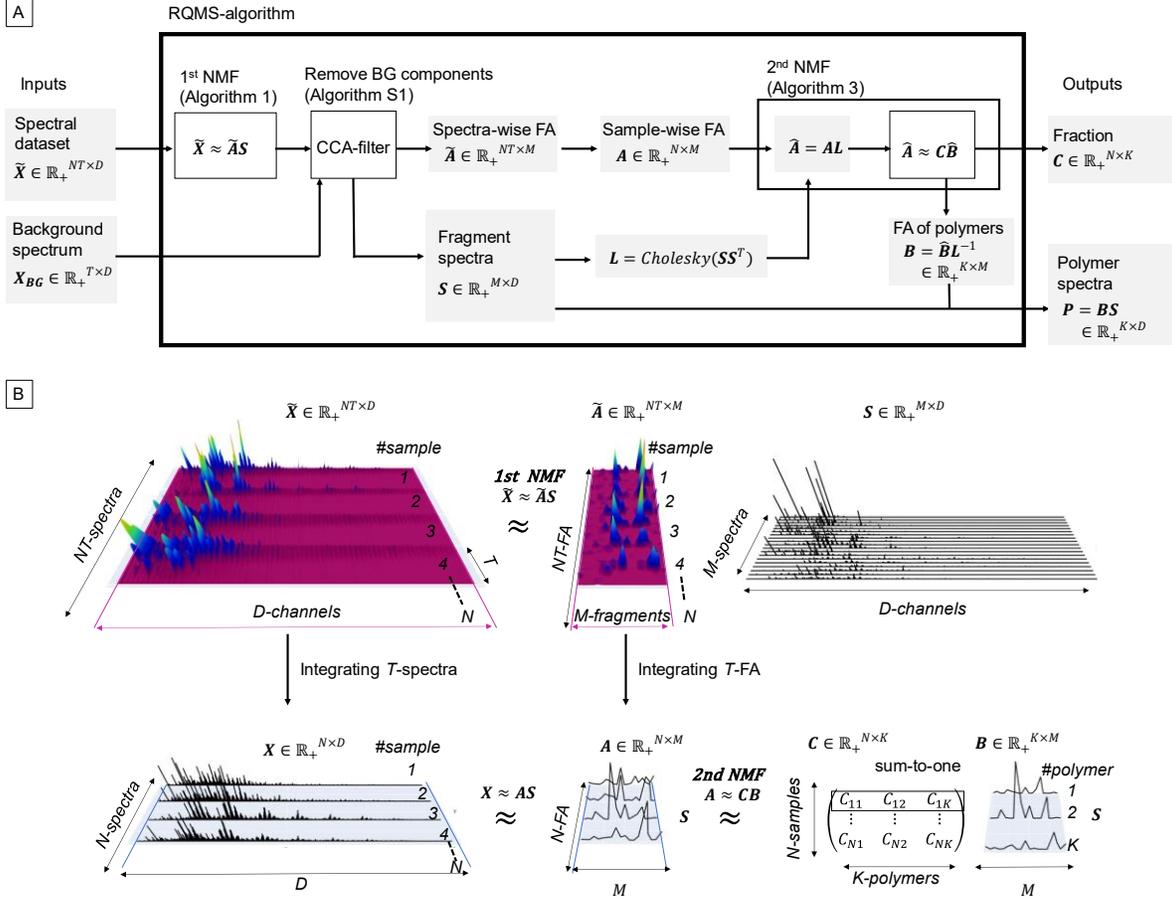


Fig. S4.

RQMS algorithm. **(A)** Flow chart, outputting C and P based on the input $\tilde{X} \in \mathbb{R}_+^{NT \times D}$. In the main text, for simplicity, we consistently assume 1D-spectral dataset $X \in \mathbb{R}_+^{N \times D}$ as the input. However, a single spectrum of pyrolysis-MS consists of T -spectra recorded at different T -temperature bands (see “Spectra formatting” section), and thus, the dataset is $\tilde{X} \in \mathbb{R}_+^{NT \times D}$ rather than $X \in \mathbb{R}_+^{N \times D}$. After the first NMF, spectrum-wise FA $\tilde{A} \in \mathbb{R}_+^{NT \times M}$ is sample-wisely integrated to sample-wise FA $A \in \mathbb{R}_+^{N \times M}$, as the temperature dependency of FA is not important for compositional analysis. Background/contaminated components are

removed by our developed CCA-filter (see “Derivation of the CCA-filter” section). To conduct the second NMF in Riemann metrics considering the non-orthogonality of \mathbf{S} (also see Fig. S6), the factorization is conducted for $\hat{\mathbf{A}} \equiv \mathbf{A}\mathbf{L} \in \mathbb{R}^{N \times M}$, where $\mathbf{L} \in \mathbb{R}^{M \times M}$ is the lower matrix obtained by Cholesky decomposition of $\mathbf{G} \equiv \mathbf{S}\mathbf{S}^T$ (see “Derivation of the second NMF” section). **(B)** The relationship between $\tilde{\mathbf{X}}$ and \mathbf{X} are graphically presented.

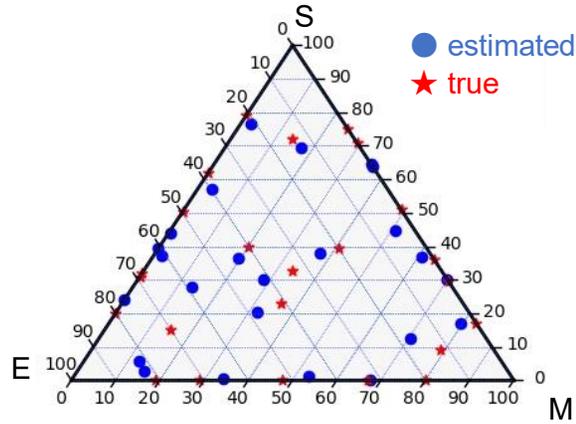


Fig. S5.

Inputting $\mathbf{X} \in \mathbb{R}_+^{N \times D}$ to RQMS algorithm derived a poor result in benchmark compositional analysis of E/M/S ternary films, suggesting the first NMF should be applied to $\tilde{\mathbf{X}} \in \mathbb{R}_+^{NT \times D}$, rather than $\mathbf{X} \in \mathbb{R}_+^{N \times D}$. All the parameters were set identical to those for Fig. 2A, except for setting temperature-band number $T = 1$ ($T = 20$ for Fig. 2A).

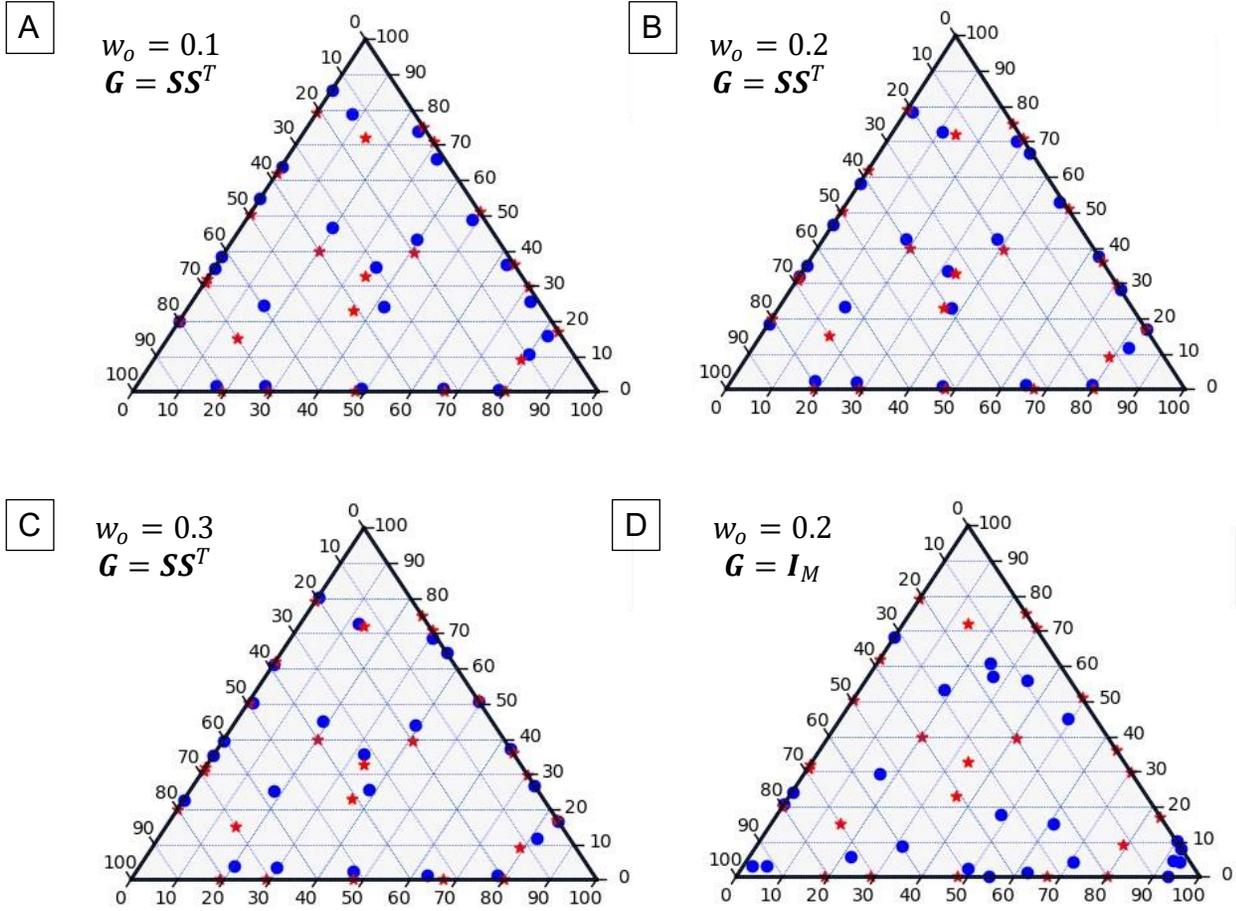


Fig. S6.

Benchmark compositional analyses of E/M/S ternary films. The residuals of the second NMF in RQMS algorithm were evaluated with Riemann metrics, i.e., $D_G(\mathbf{A}|\mathbf{CB}) = \text{Tr}[(\mathbf{A} - \mathbf{CB})\mathbf{G}(\mathbf{A} - \mathbf{CB})^T]$, where $(\mathbf{A}, \mathbf{B}, \mathbf{C})$ $\mathbf{G} = \mathbf{SS}^T$ and (\mathbf{D}) $\mathbf{G} = \mathbf{I}_M$ corresponding to Euclidean distance. $(\mathbf{A}, \mathbf{B}, \mathbf{C})$ To show the robustness of the estimated composition \mathbf{C} , the weight of orthogonal constraints in the first NMF was varied, i.e., $w_o \in \{0.1, 0.2, 0.3\}$. All the outputs were similar and less sensitive to the first NMF model selection. As the best result was obtained when $w_o = 0.2$, we consistently use this condition in this study. (\mathbf{D}) When $\mathbf{G} =$

I_M , the estimation was not reliable. Considering non-orthogonality of \mathbf{S} by introducing \mathbf{G} is thus critically important.

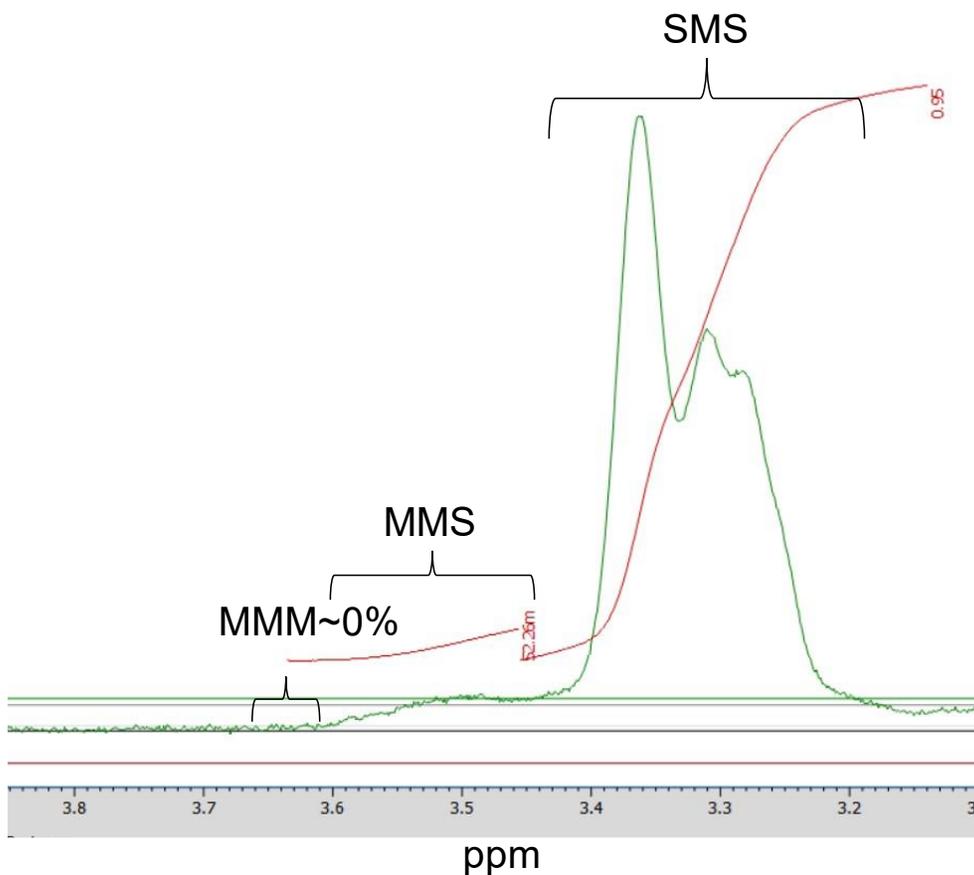


Fig. S7.

¹H NMR spectrum of commercially available (MS)_l copolymer measured in CDCl₃. The copolymer was synthesized in the presence of ethyl aluminum sesquichloride²⁶. The methoxy protons sensitive to the sequence are presented. According to the literature²⁶, MMM, MMS, SMS peaks roughly split, allowing *qualitative* analysis. As here shown, there was no MMM and very little MMS, suggesting highly alternating sequence. We did not conduct further quantitative analysis via NMR by decoupling its tacticity.

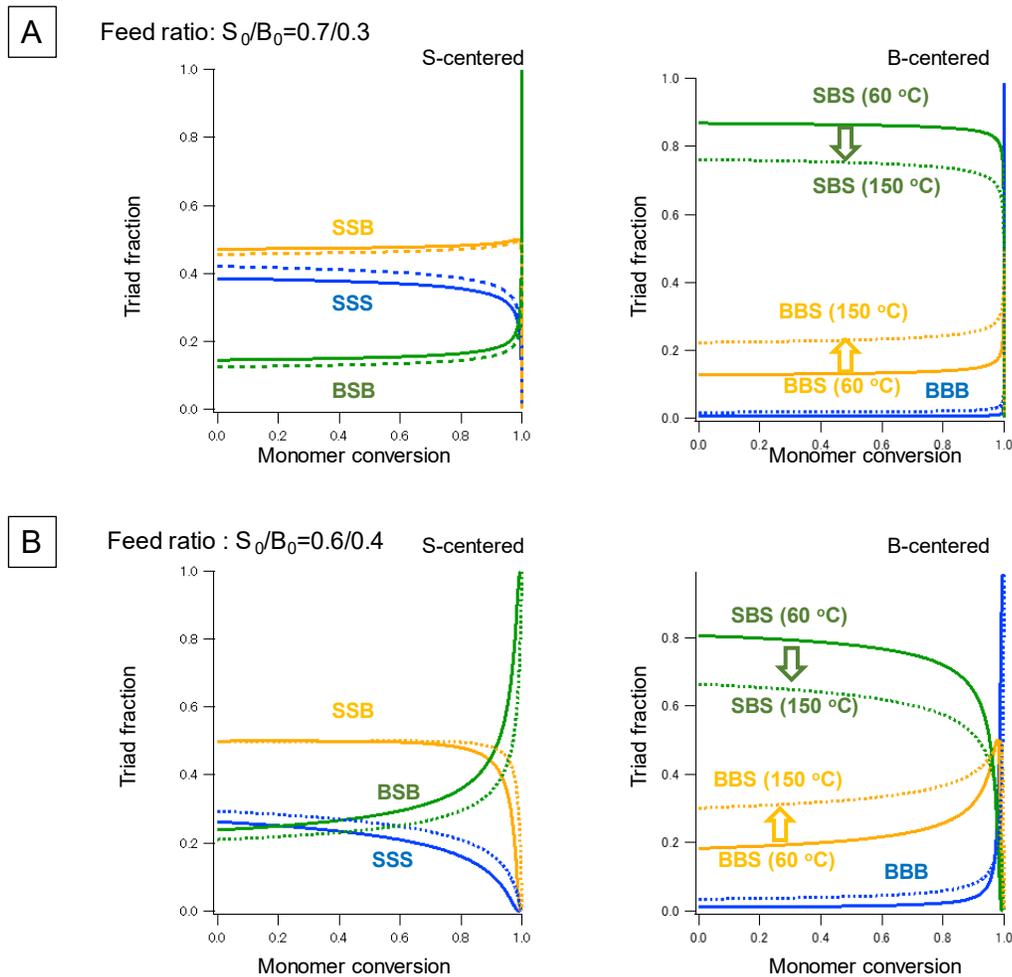


Fig. S9.

S/B instantaneously generated triad-fractions as functions of monomer conversion with the initial monomer molar fraction of **(A)** $S_0/B_0 = 0.7/0.3$ and **(B)** $S_0/B_0 = 0.6/0.4$ at polymerization temperature of 60 °C (solid lines) and 150 °C (dotted lines). S-centered and B-centered triad fractions are separately depicted. Monomer reactivity ratios are $(r_S, r_B) = (0.70, 0.17)$ and $(r_S, r_B) = (0.78, 0.34)$ at 60 °C and 150 °C polymerization temperatures²². Assume we need a copolymer having high SSSBB pentad fraction. As S/B is

inherently alternating monomer combination ($r_S < 1, r_B < 1$), achieving high SSSBB fraction is challenging. This can be intuitively understood as follows; for yielding high SSS fraction, S_0 should be sufficiently high, however, which would derive low BBS fraction owing to highly selective SBS generation. By increasing temperature from 60 °C to 150 °C, this SBS selectivity among B-centered triads is reduced, affording higher BBS generation even under high S_0 conditions as indicated by allows. This strategy was helpful for BBBSS pentad as well. In the dataset for pentad sequencing, we thus included copolymers synthesized at 150 °C with the monomer feed ratios around $S_0/B_0=0.6/0.4$ (for SSSBB) and $S_0/B_0=0.2/0.8$ (for BBBSS). For further detailed mathematical procedures, see “Prediction of sequence distribution from the monomer reactivity ratio” section. What if your monomer combination has no reported values of monomer reactivity ratio? In such case, we still can conduct triad sequencing based on a roughly prepared small dataset. Since triad fraction and monomer reactivity ratio are correlated via Alfrey-Mayo equation, the monomer reactivity ratio can be reversely calculated from the triad fraction²⁷. Based on the estimated monomer reactivity ratio, the dataset can be further designed and expanded for more advanced sequencing.

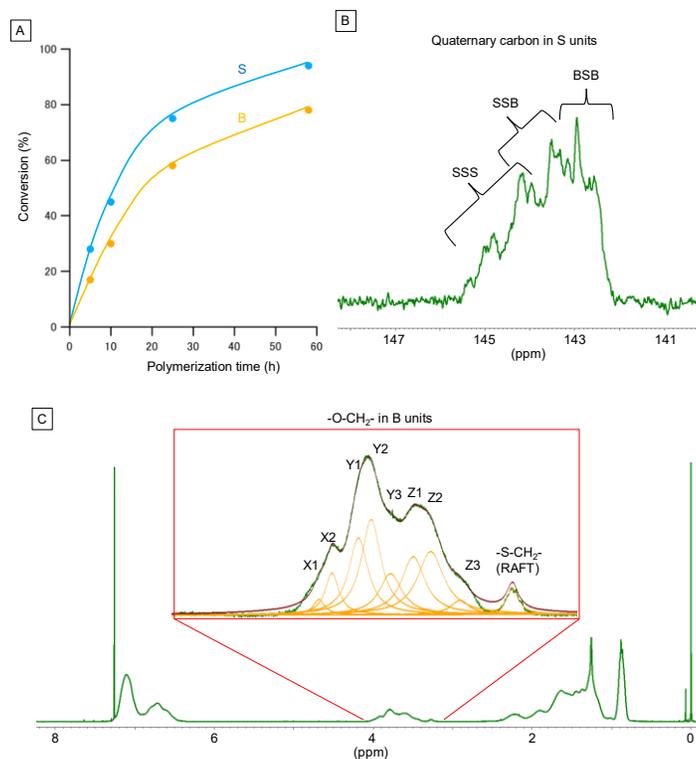


Fig. S10.

S/B random copolymer synthesis via RAFT polymerization. Polymerization conditions: $[S]_0/[B]_0/[DDMAT]_0/[AIBN]_0=20/20/0.2/0.06$ mmol in 1,4-dioxane 2 mL at 70 °C. **(A)** Monomer-polymer conversion curves of S and B. Small portions of the polymerization solution were taken out at 5 h, 10 h, 25 h, 58 h, for investigating sequence modulation along main-chains, as shown in Fig. 4B. **(B)** ¹³C NMR spectrum of S/B copolymers synthesized via RAFT copolymerization from 1/1 monomer feed ratio. Peaks of the sequence-sensitive quaternary carbon involved in S units are shown. NMR conditions: 10 wt% copolymer solution in CDCl₃ at 25 °C; 5 s relaxation time; 9000 scans in 15 h measurement time; referencing CDCl₃ central peak at 77.0 ppm. The rough peak attribution of SSS, SSB, and BSB followed the reported paper²¹. **(C)** ¹H NMR spectrum of S/B copolymers synthesized

via RAFT copolymerization. NMR conditions: 30 s of relaxation delays in CDCl₃ at 25 °C. The inset shows magnified peaks of sequence-sensitive -O-CH₂- protons involved in B monomers. The small peak on the right is attributable to -S-CH₂- protons involved in DDMAT, which is not related to sequence analysis. The overlapped peaks X, Y, Z were deconvoluted via Lorentzian peak fitting, yielding peak-area fraction of (x, y, z), which was further converted to B-centered triad fractions as described in “Determination of the S/B triad fraction via ¹H NMR” section. As simple 3-peaks fitting gave huge fitting errors, we here used 8 peaks, {X1, X2, Y1, Y2, Y3, Z1, Z2, Z3}, without a good scientific excuse. The green and brown are observed and fitted curves, respectively.

Table S1.

Hyperparameters used in this study.

Dataset	Sample number N	Temp. range (°C) Mass range(m/z)	First NMF				Second NMF		
			w_0	Merging threshold	Initial M	iteration	K	$\alpha = \beta$	p
Data S1 (Fig. 1A, S/B triad)	21	[200, 450] [50, 500]	0.2	0.99	30	5000	5	0.05	1
Data S2 (Fig. 2A, benchmark)	24	[50, 450] [50, 1000]	0.1/0.2/0.3	0.99	30	3000	3	0.1	1/1.5
Data S3 (Fig. 3, M/S triad)	30	[200, 450] [50, 500]	0.2	0.99	30	3000	5	0.05	1
Data S4 (Fig. S8, M/S/E triad)	84	[200, 450] [50, 500]	0.2	0.99	60	10000	13	0.05	1
Data S5 (Fig. 4, S/B pentad)	81	[200, 450] [100, 670]	0.2	0.99	60	10000	9	0.01	1

Table S2.

Pentad-to-triad transforming matrices, T_B and T_S . We used only T_B for downgrading RQMS pentad sequencing results to B-centered triad fraction to be comparable to NMR sequencing.

Sequence-defined copolymers	B-centered triad matrix, T_B			S-centered triad matrix, T_S		
	BBB	BBS	SBS	SSS	SSB	BSB
(BBBBB) _i	5	0	0	0	0	0
(BBBBS) _i	2	2	0	0	0	1
(BBBSS) _i	1	2	0	0	2	0
(BBSBS) _i	0	2	1	0	0	2
(BS) _i	0	0	2.5	0	0	2.5
(SSBSB) _i	0	0	2	0	2	1
(SSSBB) _i	0	2	0	1	2	0
(SSSSB) _i	0	0	1	2	2	0
(SSSSS) _i	0	0	0	5	0	0

Data S1-S5.

The full spectral datasets and sample information for Fig. 1A, Fig. 2A, Fig. 3, Fig. S8, and Fig. 4C. All the spectra recorded as CDF files have been converted into CSV files, and formatted as described in “Spectra Formatting” section, to be ready-to-use. All the spectral datasets are available in [DOI: 10.26434/chemrxiv-2022-mw76d-v2](https://doi.org/10.26434/chemrxiv-2022-mw76d-v2).

References

- 1 G. Moad, E. Rizzardo and S. H. Thang, *Aust. J. Chem.* 2012, **65**, 985–1076.
- 2 X. Zhuang, Z. Yang and D. Cordes, *Hum. Brain. Mapp.* 2020, **41**, 3807–3833.
- 3 M. Shiga, K. Tatsumi, S. Muto, K. Tsuda, Y. Yamamoto, T. Mori and T. Tanji, *Ultramicroscopy*, 2016, **170**, 43–59.
- 4 V. Y. F. Tan and C. Févotte, *IEEE Trans. Pattern. Anal. Mach. Intell.* 2013, **35**, 1592–1605.
- 5 M. Anderle, S. Roy, H. Lin, C. Becker and K. Joho, *Bioinformatics*, 2004, **20**, 3575–3582.
- 6 A. Cichocki and A.-H. Phan, *IEICE Trans. Fundam. Electron. Commun. Comput. Sci.* 2009, **E92-A**, 708–721.
- 7 K. Kimura, Y. Tanaka, M. Kudo, D. Phung and H. Li, *Mach. Learn. Res.* 2014, **39**, 129-141.
- 8 S. Boyd, N. Parikh, E. Chu, B. Peleato and J. Eckstein, *Found. Trends Mach. Learn.*, 2010, 3, 1–122.
- 9 X. Fu, W. K. Ma, K. Huang and N. D. Sidiropoulos, in *ICASSP*, 2016, 2534–2538.
- 10 D. C. Heinz and C.-I. Chang, *IEEE Trans. Geosci. Remote. Sens.* 2001, **39**, 529–545.
- 11 X. Fu, K. Huang, N. D. Sidiropoulos and W. K. Ma, *IEEE Signal Process Mag*, 2019, **36**, 59–80.
- 12 W. Liu and N. Zheng, *Pattern Recognit. Lett*, 2004, **25**, 893–897.
- 13 T. Yoshida, *Lect. Notes Comput. Sci. (including Subser. Lect. Notes Artif. Intell. Lect. Notes Bioinformatics)*, 2011, **6804 LNAI**, 214–219.
- 14 X. Fu, K. Huang, B. Yang, W. K. Ma and N. D. Sidiropoulos, *IEEE Trans. Signal Process.* 2016, **64**, 6254–6268.
- 15 C. Ding, T. Li, W. Peng and H. Park, *Proc. ACM SIGKDD Int. Conf. Knowl. Discov. Data Min.* 2006, 126–135.
- 16 J. Kim, Y. He and H. Park, *J. Glob. Optim.* 2014, **58**, 285–319.
- 17 J. M. P. Nascimento and J. M. B. Dias, *IEEE Trans. Geosci. Remote Sens.* 2005, **43**, 898–910.

- 18 L. Miao and H. Qi, *IEEE Trans. Geosci. Remote Sens.* 2007, **45**, 765–777.
- 19 F. R. Mayo and F. M. Lewis, *J. Am. Chem. Soc.* 1944, **66**, 1594–1601.
- 20 I. Skeist, *J. Am. Chem. Soc.* 1946, **68**, 1781–1784.
- 21 A. S. Brar and C. V. V. Satyanarayana, *Polym. J.* **24**, 879–887 (1992).
- 22 L. K. Kostanski and A. E. Hamielec, *Polymer (Guildf)*. 1992, **33**, 3706–3710.
- 23 H. Mutlu and J. F. Lutz, *Angew. Chem. Int. Ed.* 2014, **53**, 13010–13019.
- 24 Y. Hibi, M. Ouchi and M. Sawamoto, *Nat Commun*, 2016, **7**, 1–9.
- 25 C. I. Chang and Q. Du, *IEEE Trans. Geosci. Remote Sens.* 2004, **42**, 608–619.
- 26 H. Hirai, T. Tanabe and H. Koinuma, *J. Polym. Sci., Polym. Chem. Ed.* 1979, **17**, 843–857.
- 27 D. J. T. Hill, A. P. Lang, J. H. O'Donnell and P. W. O'Sullivan, *Eur. Polym. J.* 1989, **25**, 911–915.