

## How to encode a phylogenetic tree on a computer?

- Can I generate a random tree?
- Can I interpolate between trees?
- Can I estimate distance between trees?
- Can I infer a tree using a **deep neural network**?

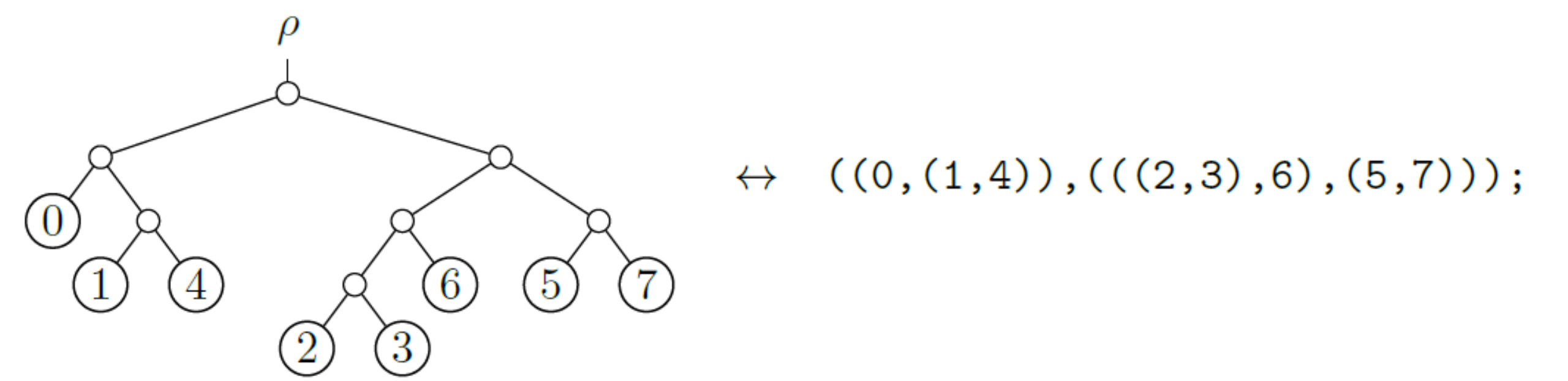


Figure 1. Current "standard" format: Newick string

## OLA encoding

- New format for encoding trees:

$$\text{OLA}_n : \left\{ \begin{array}{l} \text{rooted binary trees on} \\ 0, 1, 2, \dots, n-1 \end{array} \right\} \rightsquigarrow \left\{ \begin{array}{l} \text{vectors } (a_1, a_2, \dots) \in \mathbb{Z}^{n-1} \\ \text{with } -i < a_i < i \end{array} \right\}$$

- Encoding space is **compact**, easy to sample.
- Easy **induced distance** from Hamming distance on  $\mathbb{Z}^{n-1}$ .
- Distance in encoding space is **somewhat-low distortion** vs. SPR distance.
- Tree in Figure 1 is encoded as  $(0, -1, 2, 1, -3, -3, 5)$ .

## Related work

For all (unrooted) trees with internal & leaf nodes labeled:

- Prüfer code (1918) to vectors in  $\mathbb{Z}^{n-2}$  with  $0 < a_i \leq n$ .

For rooted binary trees with leaves labeled:

- Diaconis and Holmes (2002) define bijection to perfect matchings on leaf set.
- Chauve, Colijn, and L. Zhang (2024) encoding by certain permutations of multiset  $\{1, 1, 2, 2, \dots, n, n\}$ .
- Penn, Scheidwasser et al. (2024) "Phylo2vec" encoding by integer vectors, quadratic time complexity.

## Example encoding: Ordered leaf attachment

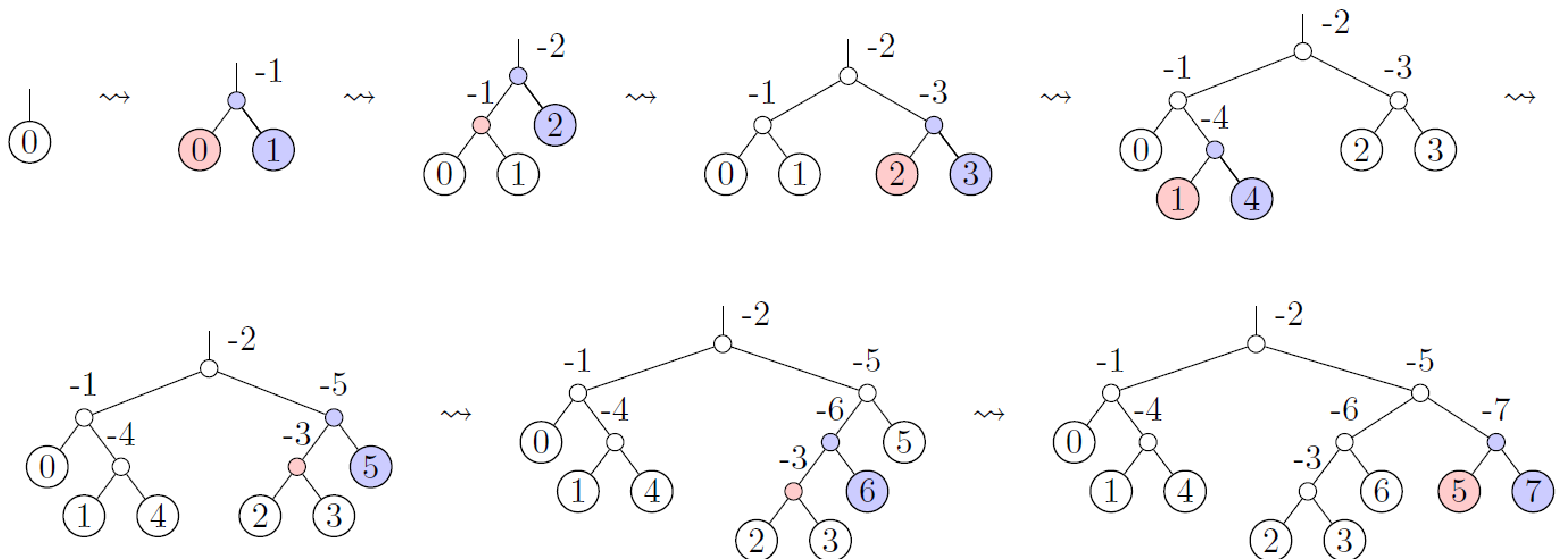


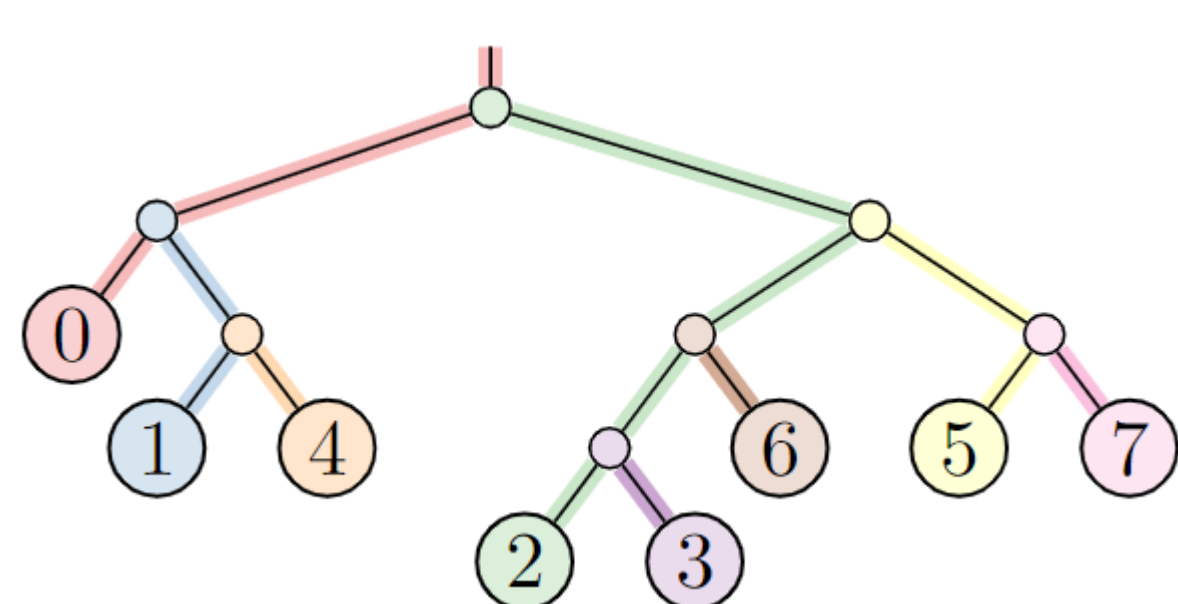
Figure 2. Constructing the tree with OLA code  $(0, -1, 2, 1, -3, -3, 5)$ .

## Results

### Theorem (Matsen–R–Zhang, 2024)

OLA encoding and decoding takes **linear time**.

- Practical implementation uses "branch decomposition" for internal node labels



### Theorem (Matsen–R–Zhang, 2024)

- If  $T'$  is a random NNI neighbor of  $T$ , then  $\mathbb{E}(d_{\text{OLA}}(T, T')) = O(1)$ .
- If  $T'$  is a random SPR neighbor of  $T$ , then  $\mathbb{E}(d_{\text{OLA}}(T, T')) = O(\text{height}(T))$ .

## OLA distance on trees

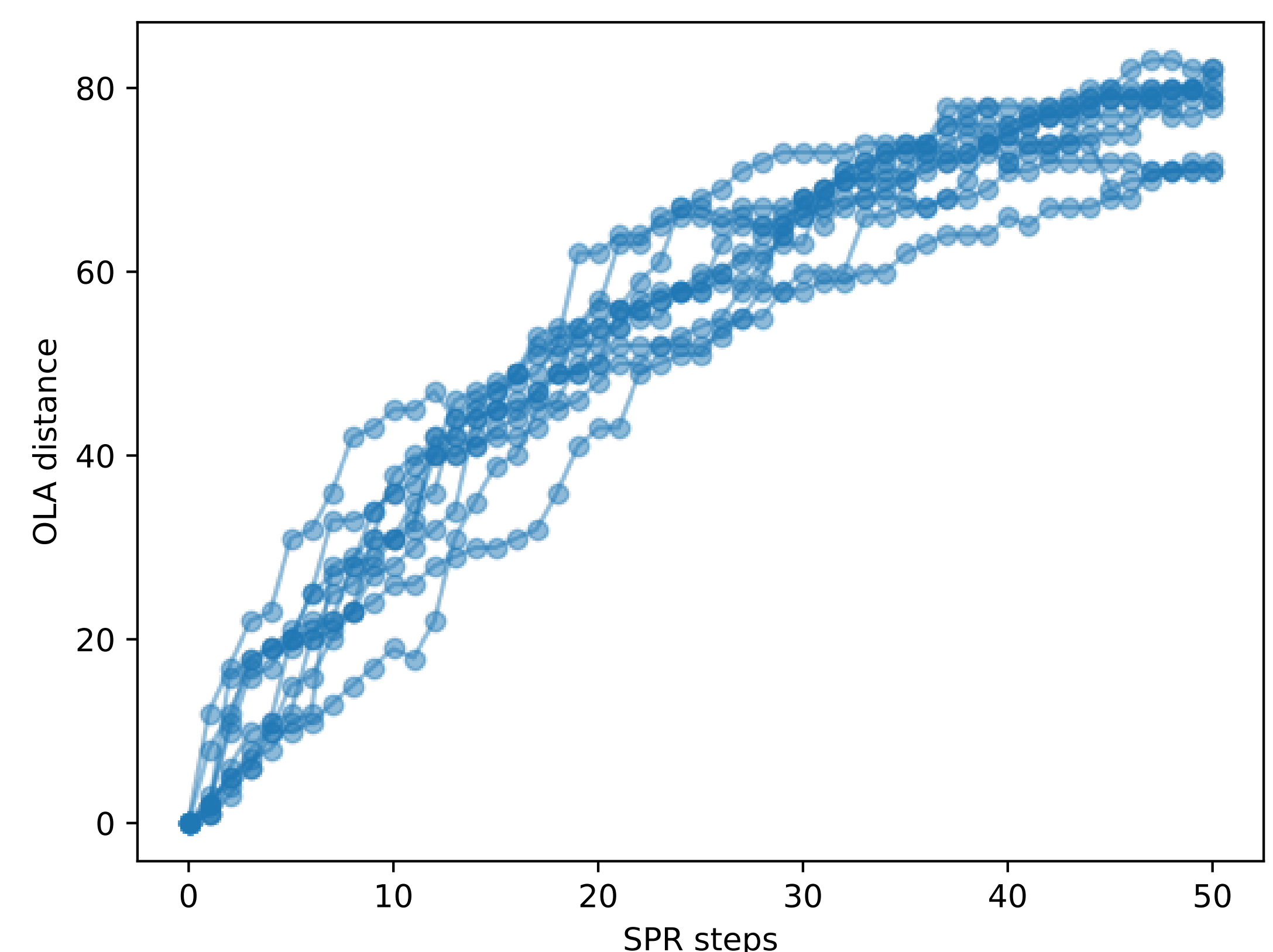


Figure 3. Experimental data from trees on 100 leaves.