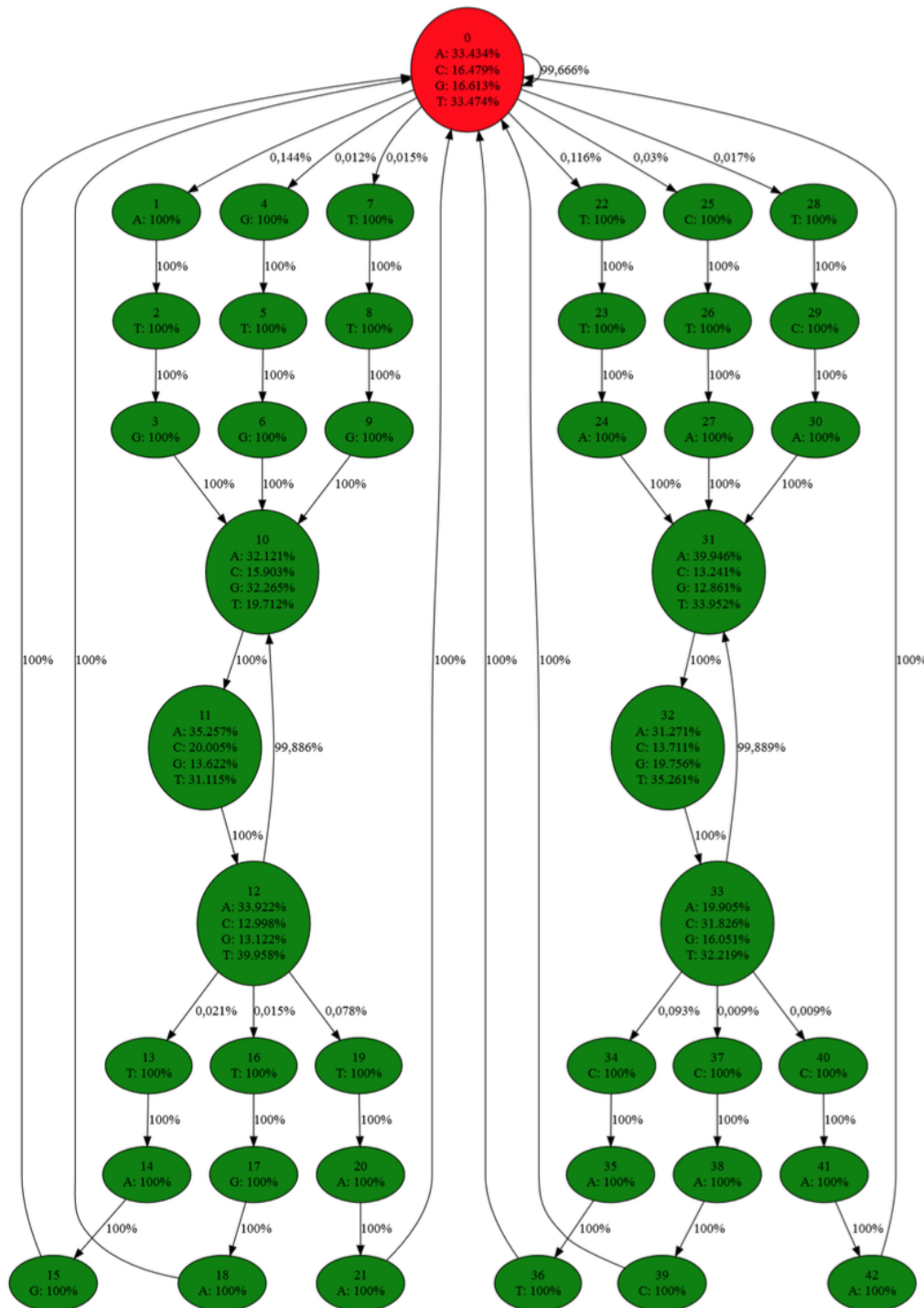
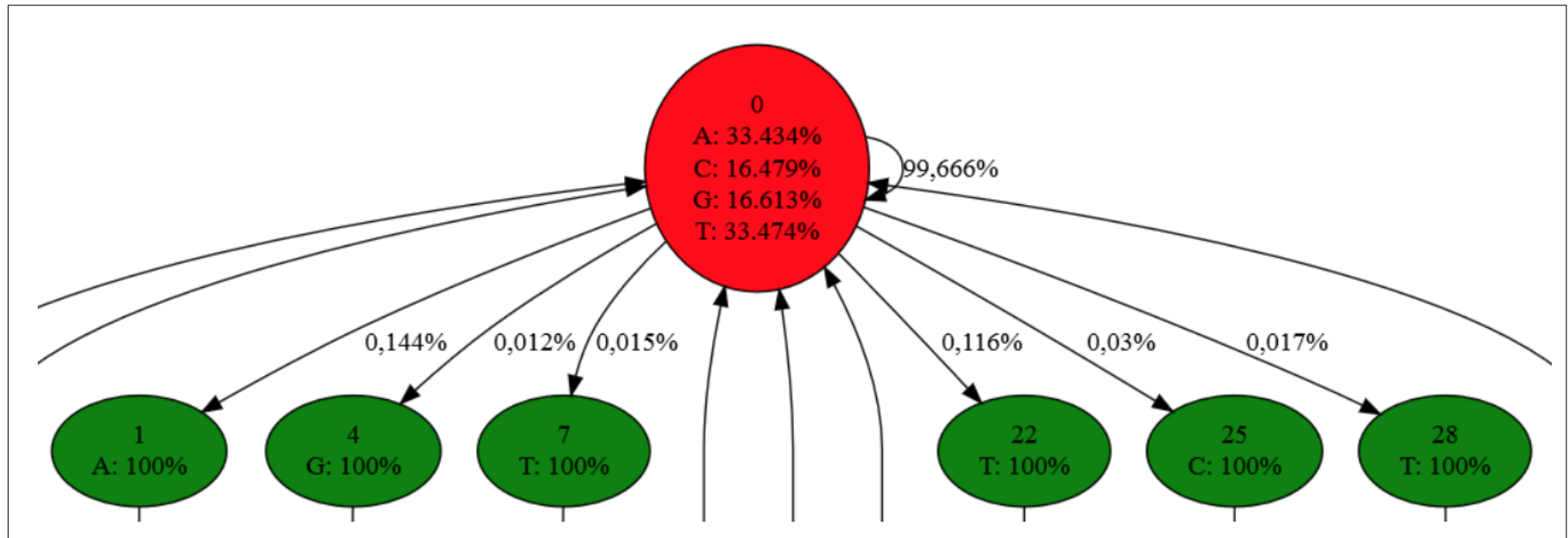


A typical model



Training by counting



In theory, we are given (\mathbf{X}, \mathbf{Z}) pairs, but we are given $(\mathbf{X}, \text{"Z"})$ pairs, where the NCR-annotations have to be translated to \mathbf{Z} s.

Also, we may ignore rare start and stop codons, i.e. a gene that starts (or ends) with a ignored start (or stop) codon does not correspond to a path in your model.

Training by counting – Typical solution

To set the transition probabilities:

$N \rightarrow N$	$N \rightarrow CCC$, where CCC is ATG	$N \rightarrow RRR$, where RRR is TTA
	$N \rightarrow CCC$, where CCC is GTG	$N \rightarrow RRR$, where RRR is CTA
	$N \rightarrow CCC$, where CCC is TTG	$N \rightarrow RRR$, where RRR is TCA

We count:

$\#(N \rightarrow N)$ = no. of occurrences of “NN” our annotations.

$\#(N \rightarrow CCC, \text{ where CCC is XYZ})$ = no. of occurrence “NCCC” in our annotations, where CCC is an annotation of XYZ (in our training data).

$\#(N \rightarrow RRR, \text{ where RRR is XYZ})$ = no. of occurrence “NRRR” in our annotations where, RRR is an annotation of XYZ (in our training data).

We compute:

$\text{Total} = \#(N \rightarrow N) + \#(N \rightarrow CCC \text{ where CCC is XYZ}) + \#(N \rightarrow RRR, \text{ where RRR is XYZ})$

We set:

$P(N \rightarrow X) = \#(N \rightarrow X) / \text{Total}$ for each of the 7 transitions