

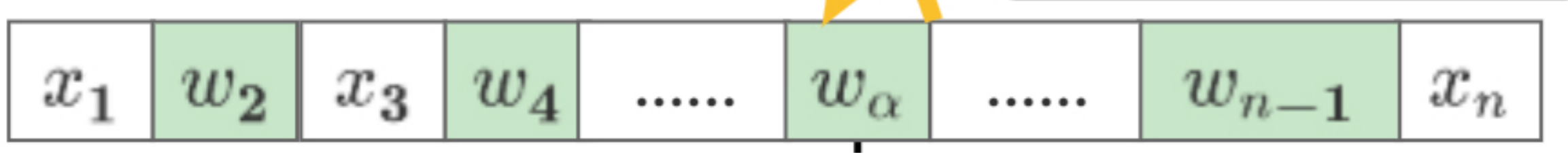
### Input Text

$$\mathbb{X} = [x_1, x_2, x_3, x_4, \dots, x_\alpha, \dots, x_{n-1}, x_n] \quad f(\mathbb{X}) = (\mathbb{Y})$$

$n$  is the index in the sentence,  $\mathbb{X}$  is the original text  
 $f$  is the classifier,  $\mathbb{Y}$  is the correct label for  $f(\mathbb{X})$

Synonyms Replacement

replacements reduction



search space reduction

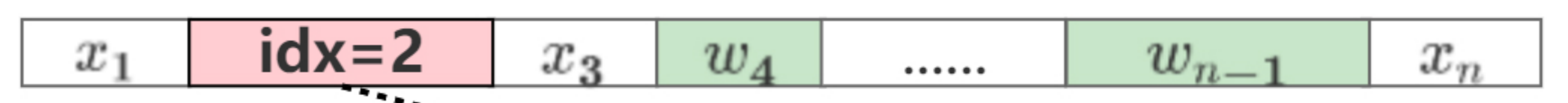
$$\mathbb{X}' = [x_1, w_2, x_3, w_4, \dots, w_{n-1}, x_n] \quad f(\mathbb{X}') \neq (\mathbb{Y})$$

$\mathbb{X}'$  is the initialization adversarial example

Replace with original input, to find important words

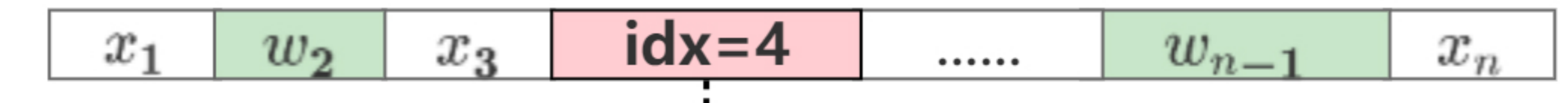
### Initialize Adversarial Example

### Synonyms Sets



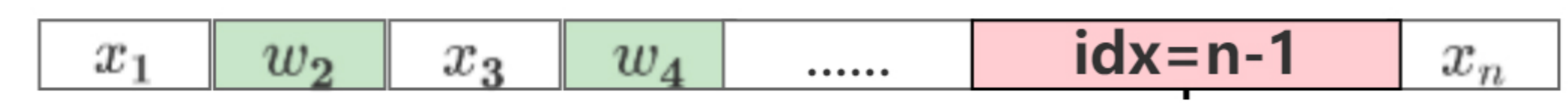
$g_1$	$w_2^1$
$g_2$	$w_2^2$
$g_3$	$w_2^3$
	$\vdots$
	$w_2^i$

Synonyms with high semantic similarity as individuals



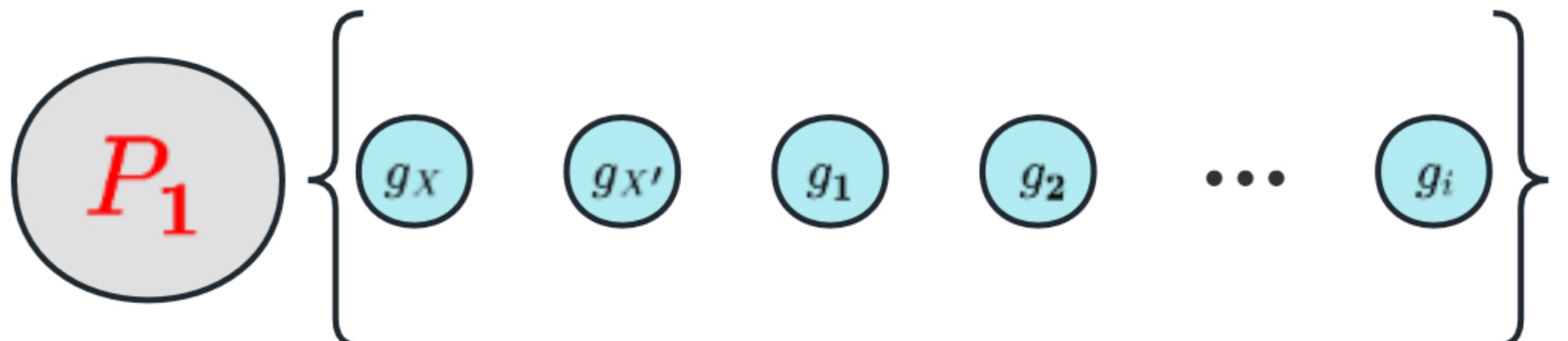
$g_4$	$w_4^1$
$g_5$	$w_4^2$
$g_6$	$w_4^3$

$\vdots$

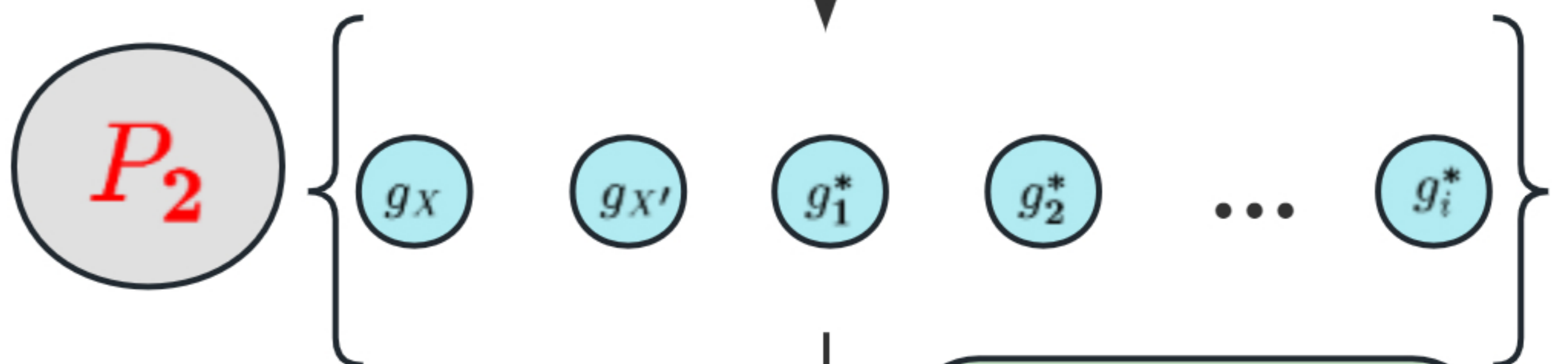


$g_i$	$w_{n-1}^1$
-------	-------------

### Generate Populations



Update



Combinatorial Optimization

Optimal Adversarial Example  $\mathbb{X}_{ADV}^*$

### Combinatorial Optimization