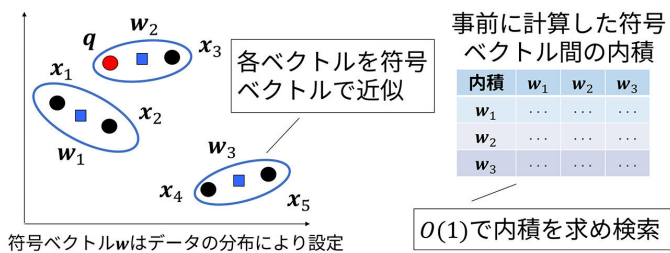


<p><b>どんな研究</b></p>	<p>ScaNNはベクトルを符号語に置き換えて量子化し、内積による類似検索を高速・高精度で行う、機械学習における重要な技術です。一方でその置き換え処理自体に時間がかかるという課題がありました。本展示ではScaNNの精度を犠牲にせずに<b>ベクトル量子化を高速に行う技術を紹介</b>します。</p>
<p><b>どこが凄い</b></p>	<p>従来のScaNNはすべてのベクトルと符号語との誤差を計算して最適な符号語に置き換える必要がありました。本技術は量子化誤差の上限に着目し、それを高速に計算することで最適になりえない候補を枝刈りします。その結果、<b>ベクトル量子化の大幅な高速化を実現</b>しました。</p>
<p><b>めざす未来</b></p>	<p>本技術は計算結果の精度を落とさずにベクトル量子化を高速化できます。内積計算は画像検索、深層学習、自然言語処理など幅広い分野で不可欠です。本技術は<b>ベクトル量子化を大規模データに対しても適用可能にし、情報処理基盤を支える中核技術</b>となることが期待されます。</p>

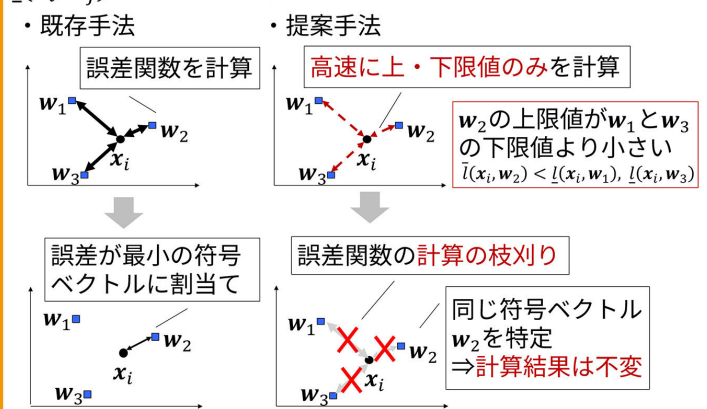
**背景：最大内積検索**

ベクトル $q$ との内積が最大になるベクトル $x$ を検索  
 ○ 画像検索、深層学習、自然言語処理などに適用可能  
 ✕ データベース内の全てのベクトルに対して内積計算が必要なため、計算コストが高くなる  
 ベクトル量子化: ベクトル $x$ を符号ベクトル $w$ に近似し検索  
 ○ 高速な検索が可能だが、✕ 検索精度が低下



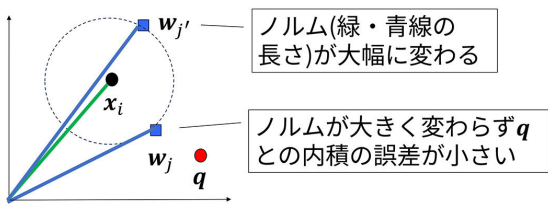
**提案手法：枝刈りによる高速化**

高速に計算できる誤差関数の上限値 $\bar{l}(x_i, w_j)$ と下限値 $\underline{l}(x_i, w_j)$ を用いて**誤差関数の計算を枝刈り**



**既存のベクトル量子化手法：ScaNN**

アイデア: ベクトルの大きさ(ノルム)が変わらないように近似し、内積の誤差を小さくすることで検索精度を向上

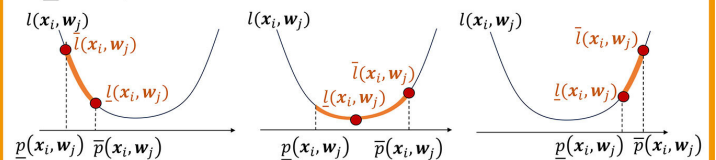


手法: 誤差関数 $l(x_i, w_j)$ が最小の符号ベクトルに割当  
 ○ 誤差関数が小さい ⇔ ノルムが大きく変わらない  
 ✕ 検索精度は向上するが、誤差関数の計算コストが高い

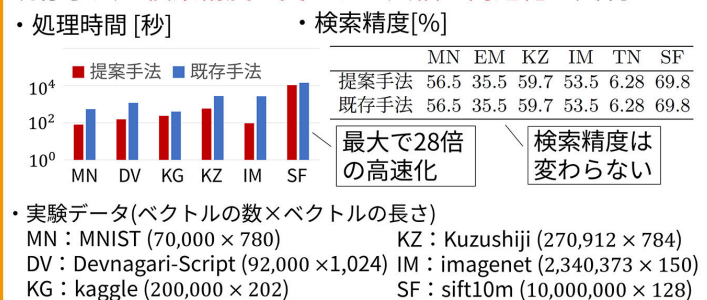
$$l(x_i, w_j) = \frac{h_i - 1}{\|x_i\|^2} (p(x_i, w_j))^2 - 2h_i p(x_i, w_j) + h_i \|x_i\|^2 + \|w_j\|^2$$

$h_i$ : 1より大きな値,  $p(x_i, w_j)$ : ベクトル $x_i$ と $w_j$ の内積

誤差関数 $l(x_i, w_j)$ は内積 $p(x_i, w_j)$ の2次関数  
 ⇒  $\bar{l}(x_i, w_j)$ と $\underline{l}(x_i, w_j)$ は、内積の上限値 $\bar{p}(x_i, w_j)$ と下限値 $\underline{p}(x_i, w_j)$ から計算



既存手法の**検索精度を変えずに大幅な高速化を実現**



関連文献

[1] Y. Fujiwara, Á. López, Y. Ida, A. Kumagai, M. Nakano, M. Nakatsuji, A. Kimura, "Fast Vector Quantization Algorithm for ScaNN", in *Proc KDD*, 2026.

連絡先

藤原 靖宏 (Yasuhiro Fujiwara) メディア情報研究部 メディア認識研究グループ