# A Deep Ranking Model for Spatio-Temporal Highlight Detection from a 360º Video

**(a) Input data processing**

360° video $X = \{x_i\}_{12}^{1}$

Input : $x_i$

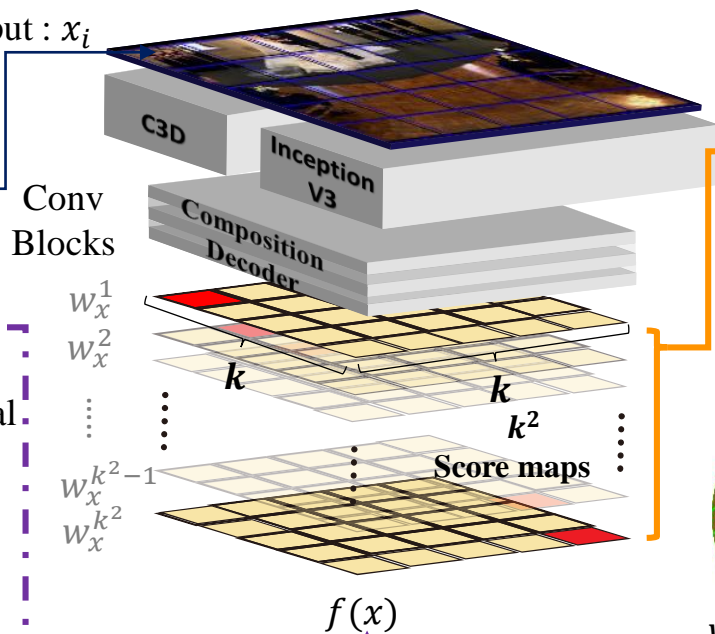Conv Blocks

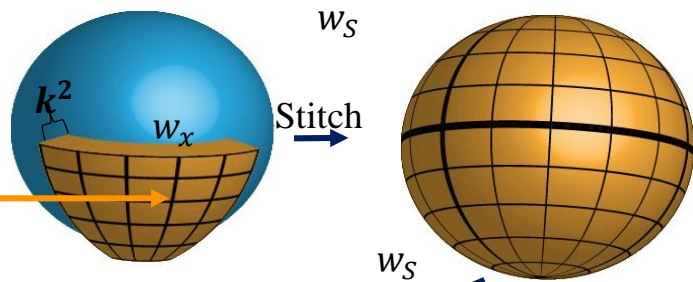C3D

Inception V3

Composition Decoder

$w_x^1$
$w_x^2$

$k$  $k$

$k^2$

Score maps

$w_x^{k^2-1}$
$w_x^{k^2}$

$f(x)$

Professional

Casual

Random

Video triplets

Deep Triplet Training

**(b) Deep ranking model**

**(c) Stitch score maps $w_X$ to a spherical score map**

$w_S$

$k^2$

$w_x$

Stitch

$w_S$

**(d) Find the best fitted region on spherical map**

composition scoring on sliding window

$w$

$w_S$

$\phi_{scale} \in \{60 \cdots, 110\}$

Multi-scale Sliding Window

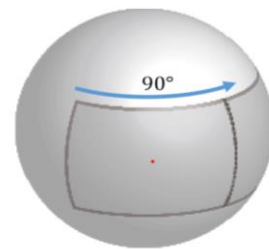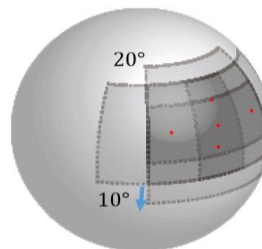Red : Best scored region
Blue : Bad scored region

**(1) Fully convolutional CVS** generates a layered spherical score maps.

**(2) Position-wise composition score** learns fidelity of views and determines which view is suitable for highlight.

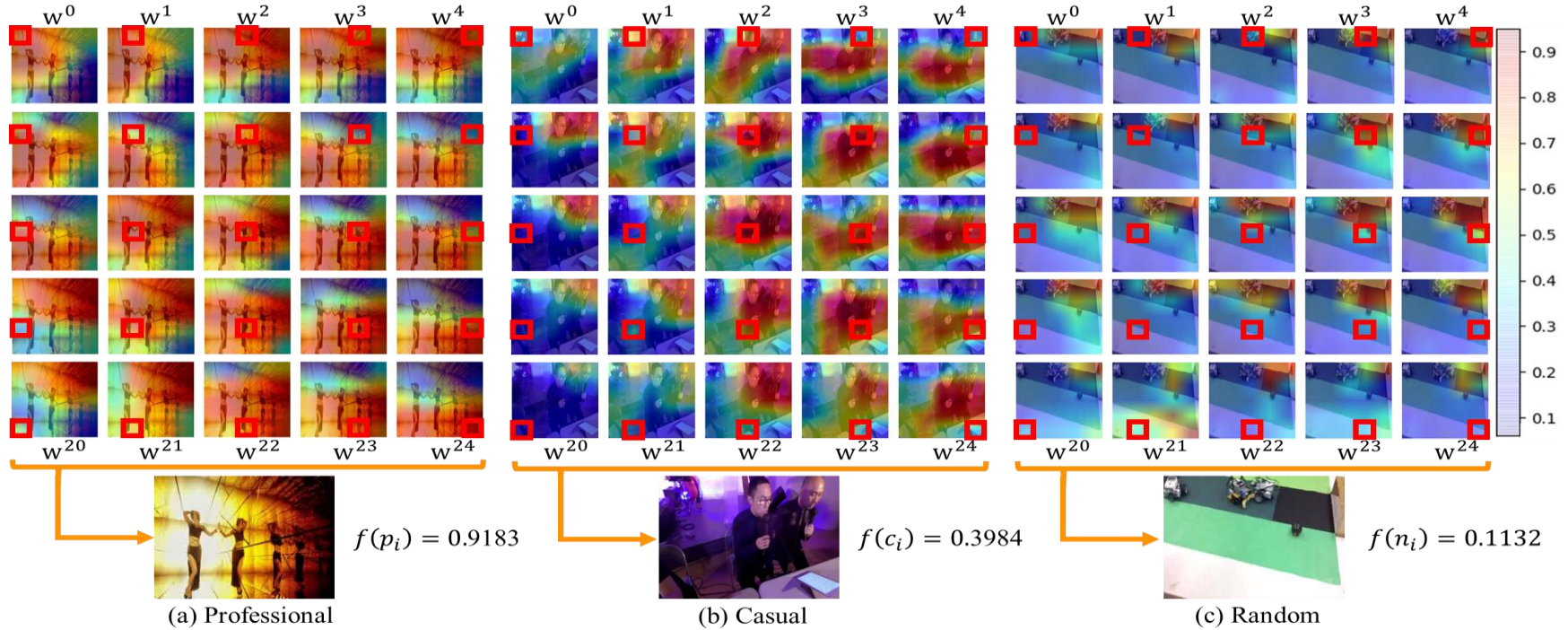**(3) Reduces a bottleneck** of normal field-of-view projection
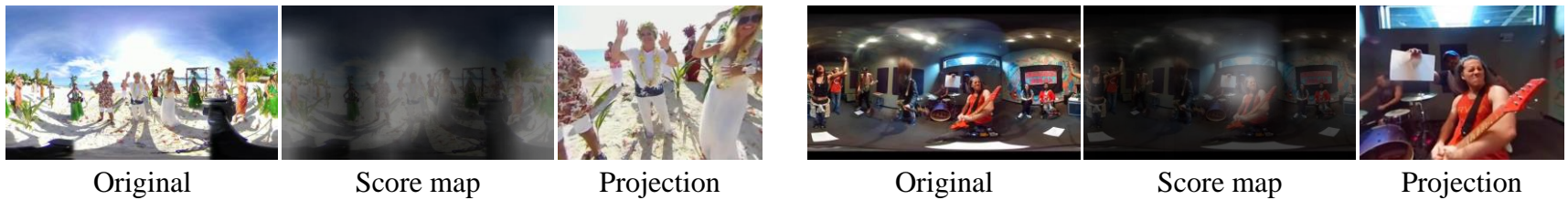
AutoCam (Su et al. 2016)

20°

10°

Ours (CVS)

90°

**# of glimpse**   198 patches   >>   **12 patches**

# Composition Score Map Learned by Triplet Deep Ranking



$f(p_i) \succ f(c_i) \succ f(n_i), \quad \forall (p_i, c_i, n_i) \in \mathcal{D}$

$w^0 \quad w^1 \quad w^2 \quad w^3 \quad w^4$

$w^{20} \quad w^{21} \quad w^{22} \quad w^{23} \quad w^{24}$

$f(p_i) = 0.9183$

$f(c_i) = 0.3984$

$f(n_i) = 0.1132$

(a) Professional

(b) Casual

(c) Random

Original  Score map  Projection

Original  Score map  Projection

**Youngjae Yu**  **Sangho Lee**  **Joonil Na**  **Jaeyun Kang**  **Gunhee Kim**

Association for the Advancement of Artificial Intelligence

SEOUL NATIONAL UNIV. VISION & LEARNING