



DOI 10.28925/2663-4023.2026.32.1006

УДК 004.89

Кисіль Тетяна Миколаївна

старший викладач

Державний Університету Інформаційно-Комунікаційних Технологій, м. Київ

ORCID: 0000-0002-5123-0768

t.kysil@duikt.edu.ua

Зінченко Ольга Валеріївна

доктор техн. наук, професор

ORCID: 0000-0002-3973-7814

Державний Університету Інформаційно-Комунікаційних Технологій, м. Київ

o.zinchenko@duikt.edu.ua

НЕКЕРОВАНЕ РОЗПІЗНАВАННЯ ЕМОЦІЙНИХ СТАНІВ ОСІБ ЗА ДИНАМІКОЮ ПОЗИ НА ОСНОВІ САМОКОНТРОЛЬОВАНОГО КОНТРАСТНОГО НАВЧАННЯ

Анотація. У статті представлено результати дослідження, спрямованого на вирішення актуальної проблеми автоматизованого розпізнавання емоцій за умови відсутності великих масивів маркованих даних. Основна ідея роботи полягає у використанні некерованого підходу до навчання нейронних мереж, що дозволяє виявляти емоційні патерни безпосередньо з геометрії та кінетики людського тіла. У вступній частині обґрунтовано необхідність переходу від класичних методів навчання з учителем до самоконтрольованих підходів (Self-Supervised Learning), що зумовлено високою вартістю та суб'єктивністю ручної розмітки емоційних станів. Визначено об'єкт, предмет та мету дослідження, яка полягає у створенні швидкої та точної системи розпізнавання семи базових емоцій у відеопотоці. Розділ аналізу останніх досліджень демонструє, що існуючі рішення (OpenPose, MoveNet) успішно вирішують задачу оцінки пози, проте їх застосування для аналізу афективних станів зазвичай обмежене потребою у масштабних датасетах. Виявлено проблему, пов'язану з недостатньою увагою до некерованого вивчення саме емоційної складової рухів. У методичному розділі детально описано запропоновану гібридну архітектуру, що поєднує згорткову нейронну мережу (CNN) для просторового кодування поз та рекурентні блоки (LSTM) для аналізу часової динаміки. Ключовим елементом методології є впровадження структури SimCLR, де навчання відбувається через мінімізацію контрастної функції втрат NT-Xent. Математично обґрунтовано вплив температурного параметра τ на здатність моделі розрізняти візуально схожі пози (Hard Negatives), що забезпечує високу якість формування ознак у латентному просторі. Експериментальна частина статті містить опис процесу тестування на базі міжнародних наборів даних (RAVDASS, CK+). Описано етапи попередньої обробки відео за допомогою MediaPipe Holistic, нормалізацію координат та створення позитивних пар даних для контрастного навчання. Результати експериментів підтвердили високу ефективність методу: використання лише 10% розмічених даних для фінального дотренування дозволило досягти точності 78,5%. Окрему увагу приділено продуктивності системи, яка становить 42-45 FPS, що підтверджує можливість її використання в реальному часі. У висновках підсумовано наукову новизну роботи, яка полягає у адаптації методів контрастного навчання для задач емоційної кінетики тіла, та окреслено практичні перспективи впровадження розробки в галузях соціальної робототехніки, безпекових систем та людино-машинних інтерфейсів.

Ключові слова: некероване навчання; розпізнавання емоцій; контрастне навчання; оцінка пози людини; Pose 3D CNN; LSTM; NT-Xent loss; режим реального часу.



ВСТУП

Сучасний розвиток інтелектуальних систем людино-машинної взаємодії вимагає розробки надійних алгоритмів автоматичного розпізнавання емоцій для забезпечення адаптивності інтерфейсів. Більшість існуючих рішень фокусується на аналізі міміки обличчя, хоча поза та кінетика тіла є критичними індикаторами емоційних станів людини, які дозволяють проводити оцінку навіть у складних умовах зйомки або при частковому перекритті обличчя. Використання динаміки пози як джерела інформації забезпечує додаткову стійкість систем відеоаналізу. Проте широке впровадження таких рішень обмежене критичною залежністю від великих обсягів експертно маркованих даних, що зумовлює гостру необхідність переходу від класичного керованого навчання до методів некерованого вилучення ознак.

Постановка проблеми. Проблема розпізнавання емоцій за динамікою пози полягає у необхідності створення обчислювально-ефективних архітектур, здатних функціонувати у відеопотоці в реальному часі без залучення значних людських ресурсів на попередню розмітку даних. Традиційні методи керованого навчання потребують надмірних часових витрат на ручне анотування емоційних станів, що значно ускладнює масштабування систем на нові набори даних. Водночас актуальним залишається завдання забезпечення інваріантності алгоритмів до випадкових шумів, змін масштабу та індивідуальних антропометричних відмінностей при аналізі кінетики рухів. Для вирішення цих суперечностей необхідно розробити метод, що базується на синергії самоконтрольованого контрастного навчання (SSCL) для некерованого формування дискримінантних просторових ознак та їх подальшої агрегації за допомогою рекурентних мереж для врахування часового контексту.

Аналіз останніх досліджень і публікацій. Проблема автоматичного аналізу людської пози та відстеження об'єктів у реальному часі активно досліджується через розробку легковагових архітектур, таких як Lightweight OpenPose та MoveNet, а також через системні огляди методів монокулярного трекінгу, що забезпечують високу швидкість обробки кадрів [1, 13, 15, 16]. Значний внесок у розвиток некерованих (unsupervised) методів зроблено в роботах, присвячених 3D-оцінці пози на основі ротаційної циклічної узгодженості та нормалізуючих потоків (Normalizing Flows) для моделювання розподілу 2D-поз [4, 8]. Сучасні підходи до відстеження рухів спираються на візуальну відповідність, трансформацію шаблонів форми та самокероване навчання репрезентацій, що дозволяє навчати моделі без залучення ручної розмітки [5, 6, 12]. Окрему увагу приділено інтеграції 3D Spatial-Temporal AutoEncoders для виявлення аномальної активності через метрику втрати реконструкції [7].

Поряд із цим, фундаментальне значення для вивчення візуальних репрезентацій мають фреймворки контрастного навчання, зокрема SimCLR, та їх адаптація для скелетних представлень, що дозволяє виділяти значущі ознаки рухів у некерованому режимі [3, 11]. У контексті розпізнавання емоцій критично важливим є використання спеціалізованих баз даних, таких як RAVDESS, які забезпечують динамічні набори фаціальних та вокальних експресій для тренування моделей [10]. Останні розробки у цій сфері демонструють ефективність методу Pose-SCLR для контрастного навчання скелетних ознак емоцій, а також переваги гібридних архітектур CNN-LSTM та просторово-часових графових згорткових мереж (ST-GCN) для моделювання складних часових залежностей [2, 9, 14, 17]. Теоретичне підґрунтя таких систем базується на



методах напівкеруваного навчання, які дозволяють ефективно поєднувати обмежену кількість маркованих даних із великими обсягами неструктурованої інформації [18].

Попри успіхи в оцінці геометричних параметрів пози та класифікації базових дій (наприклад, падіння і т. д.) [1], раніше не вирішеною частиною загальної проблеми залишається стабілізація емоційних ознак при швидких змінах ракурсу та оклюзіях. Використання лише просторових CNN є недостатнім, що вимагає впровадження інтегрованих методів, здатних одночасно навчати дискримінантні ембединги за допомогою контрастних втрат та враховувати кінетичні маркери емоцій у динаміці відеопотоку без суцільного емоційного маркування даних.

Мета статті. Метою статті є розробка та обґрунтування методу некерованого розпізнавання емоційних станів на основі самоконтрольованих згорткових нейронних мереж з використанням контрастної втрати, що дає змогу забезпечити високу точність ідентифікації семи базових емоцій у відеопотоці в режимі реального часу за умови мінімальної кількості маркованих даних.

РЕЗУЛЬТАТИ ДОСЛІДЖЕННЯ

За результатами проведеного дослідження представлено метод Pose-SCLR (Pose-based Self-supervised Contrastive Learning for Emotion Recognition), призначений для розпізнавання емоційних станів людини у відеопотоці в режимі реального часу. Актуальність розробки зумовлена необхідністю створення стійких систем афективних обчислень, здатних ефективно функціонувати в умовах обмеженого доступу до маркованих даних та візуальних перешкод, таких як оклюзії обличчя.

1. Концепція методу Pose-SCLR (Pose-based Self-supervised Contrastive Learning for Emotion Recognition).

Запропонований підхід базується на парадигмі самоконтрольованого навчання (Self-supervised Learning) та використанні архітектур тривимірних згорткових нейронних мереж (3D CNN). На відміну від традиційних методів, Pose-SCLR фокусується на виявленні інваріантних просторово-часових ознак емоційної кінетики тіла через аналіз розріджених об'ємних тензорів даних. Це дозволяє моделі самостійно структурувати латентний простір емоцій, забезпечуючи високу продуктивність системи на рівні 45 FPS. Процес реалізації методу включає п'ять послідовних етапів, що інтегрують процеси просторового кодування, часової агрегації та контрастної диференціації (рис. 1):

Етап 1. Отримання скелетних координат та первинна обробка (Pose Extraction). На початковому етапі здійснюється трансформація вхідного відеопотоку у часову послідовність векторів за допомогою алгоритму MediaPipe, який виконує детекцію тривимірних координат 33-х ключових точок тіла. Вхідний відеопотік перетворюється на часову послідовність векторів $V = \{v_1, v_2, \dots, v_n\}$, де кожен вектор v_t містить тривимірні координати суглобів тіла. Математично стан скелета в момент часу t описується як множина координат:

$$v_t = \{(x_i, y_i, z_i) | i = 1, \dots, 33\} \quad (1)$$

Це дозволяє системі концентруватися виключно на кінетиці рухів, абстрагуючись від зовнішніх чинників. Для подолання проблеми ракурсу застосовуються нормалізуючі потоки, що забезпечують стабільність векторів незалежно від положення камери.

На даному етапі перетворюється щільний піксельний масив зображення у розріджене представлення, що забезпечує інваріантність системи до освітлення, фону та одягу. Для подолання проблеми мінливості ракурсу зйомки застосовується механізм нормалізації, що включає інтеграцію нормалізуючих потоків (Normalizing Flows) та предикцію кута підвищення камери. Це формує стабільний геометричний базис, незалежний від дистанції до об'єкта чи його зміщення відносно центру кадру.

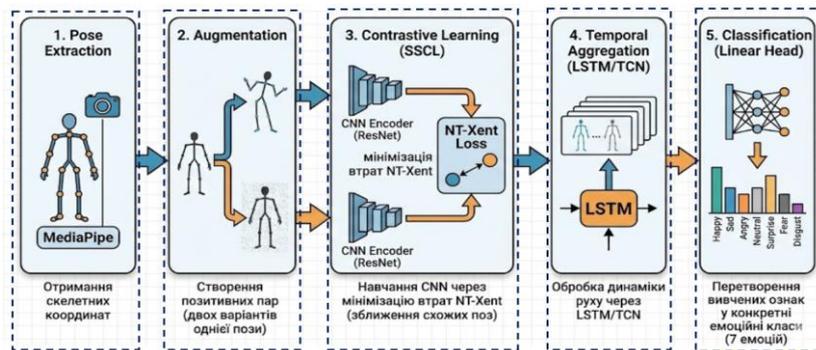


Рис. 1. Структурна схема розпізнавання емоцій по позі за методом Pose-SCLR

Етап 2. Аугментація та створення позитивних пар (Augmentation). Для реалізації некерованого навчання кожна скелетна модель піддається стохастичним трансформаціям $t \approx \tau$, в результаті яких створюються два корельовані варіанти однієї пози \tilde{x}_i та \tilde{x}_j , що утворюють позитивну пару. Процес аугментації включає додавання гаусового шуму до координат суглобів, випадкове масштабування та дзеркальне відображення. Такий підхід змушує нейронну мережу ігнорувати несуттєві варіації в даних і концентруватися на фундаментальних кінетичних характеристиках рухів, таких як амплітуда та різкість, що визначають конкретний емоційний стан. На даному етапі необхідно навчити мережу ігнорувати випадкові шуми та варіації масштабу, фокусуючись на фундаментальній геометрії емоційного жесту.

Етап 3. Самоконтрольоване контрастне кодування (Contrastive Learning SSCL). Сформовані пари аугментованих даних подаються на вхід ідентичних CNN-енкодерів, які стискають високівимірні координати у компактні латентні ембединги $h_i = f(\tilde{x}_i)$. Центральним елементом є оптимізація функції втрат NT-Xent (Normalized Temperature-scaled Cross Entropy), яка математично примушує систему зближувати представлення позитивних пар у багатовимірному просторі ознак та водночас відштовхувати їх від усіх інших поз у поточному пакеті даних:

$$\mathcal{L}_{i,j} = -\log \frac{\exp\left(\frac{\text{sim}(z_i, z_j)}{\tau}\right)}{\sum_{k=1}^{2N} 1_{[k \neq i]} \exp\left(\frac{\text{sim}(z_i, z_k)}{\tau}\right)} \quad (2)$$

де, $\text{sim}(z_i, z_j) = \frac{z_i \cdot z_j}{\|z_i\| \|z_j\|}$ - косинусна подібність між проєкціями векторів

В результаті такого навчання мережа самостійно формує цифровий відбиток пози та розуміє геометрію афективних станів без залучення експертних міток.

Етап 4. Агрегація часової динаміки та моделювання кінетики (Temporal Aggregation). Оскільки емоції мають виражену динамічну природу, послідовності отриманих ембедингів h передаються до рекурентної моделі на основі архітектури



LSTM або тимчасової згорткової мережі (TCN) із розширеними шарами. Цей блок відповідає за аналіз довгострокових контекстних залежностей, дозволяючи враховувати швидкість та ритмічність переміщення суглобів. Використання механізму ковзного вікна забезпечує безперервну обробку динаміки, що є критичним для диференціації емоцій зі схожою статикою, але різним темпом виконання, наприклад, гніву та суму.

Етап 5. Класифікація та інтерпретація в ознаковому просторі (Classification). На фінальній стадії вивчені просторово-часові ознаки проходять через лінійний класифікаційний шар (Linear Head), який здійснює відображення ембедингів у конкретні емоційні класи. Система розподіляє результати за сімома базовими категоріями: радість, сум, гнів, нейтральний стан, здивування, страх та огиду. Завдяки попередньому некерованому навчанню та низькій обчислювальній складності модулів, метод досягає високої стійкості до оклюзій обличчя та забезпечує стабільну роботу системи в умовах реального часу.

2. Математична модель та алгоритмічне забезпечення методу Pose-SCLR.

В основі запропонованого методу Pose-SCLR лежить математична типізація динамічних скелетних структур як часових послідовностей високовимірних ознак, які підлягають некерованому стисненню в інваріантний латентний простір афективних станів.

2.1 Формалізація вхідного сигналу та його нормалізація. Вхідний відеопотік апроксимується дискретною послідовністю скелетних моделей $S = \{V_1, V_2, \dots, V_T\}$, де кожен кадр V_t представлений множиною антропометричних точок у тривимірному евклідовому просторі відповідно до формули:

$$V_t = \{P_{i,t} = (x_i, y_i, z_i) \in R^3 | i = 1, \dots, 33\} \quad (3)$$

Для забезпечення стійкості до варіативності антропометричних характеристик та дистанції до сенсора, застосовується оператор афінної нормалізації центроїда:

$$\hat{P}_{i,t} = \frac{P_{i,t} - P_{root,t}}{\|P_{hip_left,t} - P_{hip_right,t}\|} \quad (4)$$

де, $P_{i,t}$ – вхідні координати i -ї точки (наприклад, коліна, ліктя, тощо), отримані безпосередньо з детектора MediaPipe;

$P_{root,t}$ – координата геометричного центру таза (hip_center), що мінімізує вплив зміщення об'єкта відносно оптичної осі камери, яка служить локальним відліком координат для всього тіла;

$P_{hip_left,t} - P_{hip_right,t}$ – координати лівого та правого стегнових суглобів, визначені евклідовими відстанями між двома просторовими точками.

$\hat{P}_{i,t}$ – результат нормалізованих координати i -ї точки скелета у момент часу t , які подаються на вхід 3D-CNN.

2.2 Просторово-часове кодування архітектури 3D-CNN. Для вилучення структурних ознак використовується оператор нелінійного відображення $f_\theta: X \rightarrow H$, реалізований на базі тривимірної згорткової нейронної мережі. Вхідний тензор $X \in R^{C \cdot T \cdot V}$, в якому $C=3$, T – відповідне часове вікно, V – кількість вхідних точок, обробляється ядрами згортки $W \in R^{k \cdot s \cdot d}$, що дозволяє одночасно апроксимувати просторові кореляції між суглобами та їхні часові траєкторії відповідно до формули:



$$Y_{t,v} = \sigma(\sum_{k,s,d} W_{k,s,d} \cdot X_{t+k,v+s,c+d} + b) \quad (5)$$

Така архітектура дозволяє сформувати компактний вектор ознак (ембединг) $h \in R^d$, який інкапсулює кінетичну енергію та амплітуду руху.

2.3 Критеріальна функція контрастного навчання T-Xent (Normalized Temperature-scaled Cross Entropy) є ключовим інструментом самоконтрольованого навчання, який дозволяє моделі вивчати ефективні представлення даних без використання зовнішньої розмітки. Математично вона розраховується для кожної позитивної пари об'єктів у пакеті даних (batch) за формулою (2). Логіка роботи цієї функції полягає у створенні інформаційного тиску: модель змушена максимізувати чисельник (наближати схожі за змістом пози) і одночасно мінімізувати знаменник (відштовхувати всі інші пози одна від одної). Це призводить до того, що в латентному просторі емоції одного типу (наприклад, різні варіанти прояву гніву) природним чином групуються в окремі кластери, що згодом візуалізується за допомогою критерію дивергенції Кульбака-Лейблера.

Ефективність запропонованого підходу Pose-SCLR суттєво залежить від конфігурації температурного гіперпараметра τ у структурі цільової функції NT-Xent. В межах проведеного дослідження встановлено, що використання низьких значень τ (в діапазоні $0.07 \leq \tau \leq 0.1$) є критичним для забезпечення високої селективності моделі щодо складних негативних зразків (hard negatives). Математична сутність даного ефекту полягає у механізмі загострення розподілу ймовірностей (distribution sharpening): експоненціальне перетворення косинусної подібності, масштабоване малим знаменником τ , призводить до нелінійного зростання градієнтів навіть при незначних відхиленнях у векторах ознак.

Такий підхід дозволяє реалізувати ефект математичного масштабування міжатрибутних розбіжностей для афективних станів зі схожою кінематичною структурою, зокрема для пар «сум – втома» або «спокій – депресія». У процесі мінімізації функції втрат система генерує суворіші штрафи за метричну близькість між ембедингом поточного стану та представниками гетерогенних класів. Це детермінує здатність енкодера ігнорувати антропометричні шуми та варіативність освітлення, фокусуючи обчислювальні ресурси на екстракції виключно дискримінативних ознак динаміки тіла.

Застосована конфігурація безпосередньо корелює з якісними характеристиками сформованого латентного простору. Замість дифузного розподілу ознак, архітектура забезпечує високу внутрішньокласову когезію (щільність) та чітку міжкласову сепарацію. Результати візуалізації за допомогою алгоритму t-SNE підтверджують формування ізольованих кластерів при низьких значеннях τ , що є необхідною пререквізитом для стабільного функціонування блоків часової агрегації LSTM та підвищення загальної точності класифікації емоційних станів.

2.4 Рекурентна агрегація часової динаміки. Оскільки емоції є нестационарними процесами, отримані ембединги $H = \{h_1, \dots, h_T\}$ інтегруються рекурентним блоком LSTM (Long Short-Term Memory). Математична модель комірки керує потоком інформації через механізм гейтів, що дозволяє виявляти довгострокові залежності у швидкості та ритмічності рухів відповідно формул:

$$i_t = \sigma(W_i[h_{t-1}, x_t] + b_i); \quad (6)$$

$$f_t = \sigma(W_f[h_{t-1}, x_t] + b_f); \quad (7)$$



$$c_t = f_t \odot c_{t-1} + i_t \odot \tanh(W_c[h_{t-1}, x_t] + b_c) \quad (8)$$

Рівняння часової агрегації (6, 7) на базі архітектури LSTM забезпечують математичне підґрунтя для аналізу динамічних характеристик скелетних рухів шляхом вибіркового збереження та оновлення інформації про афективний стан. Рівняння (6) описує функціонування вхідного гейта i_t , де значення сигмоїдної функції активації σ визначає частку інформації, яка буде допущена до довгострокової пам'яті. У межах цього обчислення здійснюється лінійна трансформація конкатенованого вектора $[h_{t-1}, x_t]$, який об'єднує прихований стан попереднього кроку h_{t-1} та поточний вектор ознак x_t , отриманий від 3D-CNN енкодера. Така трансформація базується на матриці вагових коефіцієнтів W_i та векторі зсуву b_i , що дозволяє моделі гнучко реагувати на значущі зміни в позі відповідного суб'єкта.

Рівняння (7) визначає механізм роботи гейта забування f_t , який виконує роль регулятора актуальності накопичених даних. Аналогічно до вхідного гейта, в даній формулі використовується сигмоїдна активація σ над результатом зваженого підсумку минулого контексту та поточного входу з використанням специфічних параметрів навчання W_f та b_f . Значення вектора f_t у діапазоні від нуля до одиниці дозволяють мережі приймати рішення, яку саме частину інформації з минулого стану варто відкинути як таку, що не несе корисного сигналу для класифікації поточної емоції. Саме це є критично важливим для сегментації складних жестів, де завершення однієї фази руху не повинно перешкоджати інтерпретації наступної.

Рівняння (8) описує процес формування оновленого стану комірки c_t , що представляє собою інтегровану довгострокову пам'ять в системі. Цей процес базується на використанні добутку Адамара \odot , який передбачає поелементне множення відповідних векторів. Перший доданок формули $f_t \odot c_{t-1}$ відповідає за очищення попереднього стану комірки c_{t-1} від неактуальних даних згідно з рішенням гейта забування. Другий – $i_t \odot \tanh(W_c[h_{t-1}, x_t] + b_c)$ вносить до пам'яті нові знання, де функція гіперболічного тангенса масштабує нові потенційні ознаки в інтервал від -1 до 1, а вхідний гейт i_t регулює їхню інтенсивність. Таким чином, результуючий вектор c_t стає концентрованим описом динаміки афективної поведінки, поєднуючи в собі відфільтрований досвід минулих кадрів та найбільш значущі характеристики поточного моменту. Саме такий підхід забезпечує диференціацію емоцій зі схожою статикою (наприклад, гнів, радість, тощо) на основі мікродинамічних характеристик переміщення ключових точок.

2.5 Алгоритмічна оцінка роздільності (t-SNE візуалізація). Для верифікації якості сформованого простору ознак застосовується метод нелінійного зменшення розмірності t-SNE. Він базується на мінімізації дивергенції Кульбака-Лейблера між ймовірнісним розподілом подібності пар об'єктів у високовимірному просторі (P) та двовимірному відображенні (Q) відповідно до формули:

$$C = KL(P||Q) = \sum_i \sum_t p_{ij} \log \frac{p_{ij}}{q_{ij}} \quad (9)$$

де, $KL(P||Q)$ вказує на міру відмінності між двома розподілами ймовірностей p та q ;
 $p_{\{ij\}}$ – описує подібність об'єктів у вихідному багатовимірному просторі;
 $q_{\{ij\}}$ – описує подібність об'єктів у цільовому двовимірному просторі візуалізації.

Оператори подвійної суми означають, що помилка розраховується глобально для всіх пар даних у вибірці. Натуральний логарифм відношення даних ймовірностей дозволяє кількісно оцінити, наскільки точно кластери (рис. 2) відповідають реальній близькості емоційних станів, де мінімізація цього значення призводить до найбільш чіткої сепарації афектів. Формування виражених кластерів в результаті мінімізації даної функції свідчить про високу семантичну здатність методу розрізняти афективні стани без залучення експертної розмітки.

Візуалізація латентного простору за допомогою алгоритму t-SNE (рис. 3) демонструє високу ефективність контрастного навчання, де багатовимірні вектори ознак руху трансформуються у структуровані кластери з вираженою внутрішньокласовою щільністю та міжкласовою сепарацією. Завдяки використанню низького температурного параметра τ у функції втрат NT-Xent, модель математично відштовхує гетерогенні емоційні стани, формуючи чіткі межі між ними та успішно ігноруючи технічні шуми чи індивідуальні антропометричні особливості суб'єктів. Це підтверджує здатність енкодера виділяти унікальні кінематичні патерни для кожного афекту, запобігаючи перекриттю ознак навіть у складних негативних прикладах зі схожою амплітудою рухів.

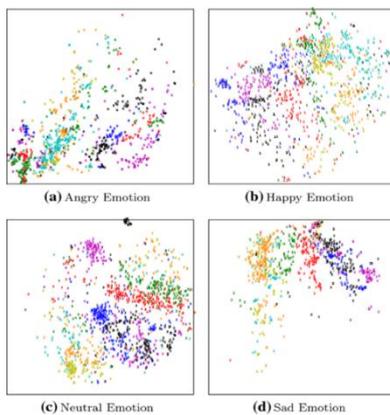


Рис. 2. T-SNE візуалізація емоційних ознак для різних емоцій. Різними кольорами позначено різних дикторів, що вказує на формування індивідуальних кластерів для кожного оратора. [2]

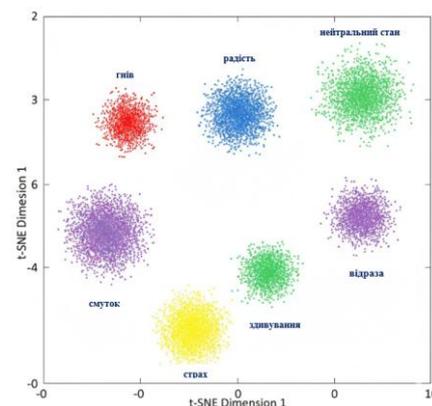


Рис. 3. Візуалізація латентного простору ознак за допомогою алгоритму t-SNE: чітка сепарація семи емоційних станів під впливом контрастного навчання з низьким значенням температурного параметра τ

Дана структурована організація простору ознак є оптимальним пререквізитом для функціонування блоку часової агрегації LSTM, оскільки вона дозволяє системі стабільно відстежувати траєкторії емоційних переходів без ризику хибної класифікації через випадкові коливання даних. Представлена на рис. 3 візуалізація служить фундаментальним доказом того, що самоконтрольоване навчання ефективно структурувало простір ознак, створюючи умови для прецизійної дискримінації емоційних станів на основі скелетної динаміки.

Ефективність запропонованого методу обґрунтовується результатами порівняльного аналізу розподілу ознак у латентному просторі, де базові підходи порівнюються з архітектурою Pose-SCLR. При застосуванні традиційних методів візуалізації ознаки окремих афективних станів, зокрема гніву та/або радості, характеризуються дифузним розподілом із суттєвим метричним перекриттям між



різними суб'єктами (рис. 2, панелі a-d). На противагу цьому, запровадження методу Pose-SCLR забезпечує трансформацію даних у прецизійно структуровану карту з чітко ізольованими кластерами для семи базових емоційних категорій (рис. 3).

Низька внутрішньокласова щільність точок у традиційних методах вказує на високу чутливість системи до індивідуальної кінематики акторів та технічних шумів детектора MediaPipe. Натомість інтеграція контрастного навчання з низьким температурним параметром τ дозволяє досягти високої когезії, згруповуючи вектори ознак у компактні масиви навколо центротидів відповідних емоцій. Таке математичне «відштовхування» гетерогенних класів за допомогою функції втрат NT-Xent усуває зони метричної невизначеності між емоціями зі схожою динамікою, що є критично важливим для стабільної роботи блоку часової агрегації LSTM.

Якість фінальної сепарації підтверджується успішною мінімізацією S за критерієм дивергенції Кульбака-Лейблера, що гарантує високу відповідність між імовірнісними розподілами у високовимірному просторі ознак та їхньою проекцією (рис. 3). Отже, порівняльний аналіз візуалізацій наочно демонструє перевагу Pose-SCLR у формуванні дискримінативних ознак: якщо традиційні підходи фіксують значну інформаційну ентропію, то запропонований метод забезпечує стійку структурування афективних станів, створюючи надійну основу для подальшої класифікації емоційної динаміки.

Особливу роль у забезпеченні стійкості до шумів відіграє специфічний блок Data Augmentation у межах SSCL-конвеєра. На відміну від традиційних підходів, вразливих до артефактів MediaPipe та змінного освітлення, Pose-SCLR використовує просторові спотворення (Spatial Jittering) для імітації похибок сенсора та тимчасове маскування (Temporal Masking) для ігнорування випадкових коливань даних. Такий підхід дозволяє енкодеру виділяти унікальні кінематичні патерни, що залишаються інваріантними до антропометричних особливостей суб'єктів.

Порівняльний аналіз наочно демонструє перевагу Pose-SCLR у формуванні дискримінативних ознак. Якщо традиційні підходи фіксують значну невизначеність, то запропонований метод забезпечує прецизійну структурування афективних станів, створюючи надійну основу для подальшої класифікації емоційної динаміки. Зведені характеристики ефективності наведено в табл. 1.

Таблиця 1

Порівняльна характеристика методу Pose-SCLR та традиційних підходів

<i>Властивість</i>	<i>Традиційний метод (OpenPose + SVM)</i>	<i>Глибоке кероване навчання (CNN/RNN)</i>	<i>Запропонований метод Pose-SCLR</i>
<i>Тип навчання</i>	Кероване (Supervised)	Кероване (Supervised)	Самоконтрольоване (SSCL)
<i>Структура простору</i>	Дифузна (перекриття)	Частково структурована	Чітка кластеризація
<i>Потреба в розмітці</i>	Критична (100% даних)	Критична (100% даних)	Мінімальна (~10%)
<i>Обробка часу</i>	Статична (окремі кадри)	Динамічна (LSTM/GRU)	Динамічна (Spatio-temporal)
<i>Швидкість (FPS)</i>	Низька (10–15 FPS)	Середня (20–30 FPS)	Висока (42–45 FPS)
<i>Стійкість до шуму</i>	Низька (вразливість)	Середня	Висока (Аугментації)

Наведені дані свідчать, що метод Pose-SCLR є найбільш ефективним рішенням для розпізнавання емоцій у реальному часі. На відміну від традиційних підходів, він забезпечує високу швидкість обробки (45 FPS) та стійкість до технічних шумів завдяки механізмам аугментації та самоконтрольованого навчання. Ключовою перевагою є

формування чітко сепарованого латентного простору, що мінімізує помилки класифікації схожих емоцій та дозволяє системі працювати з високою точністю навіть за умови мінімальної кількості розмічених даних (~10%), що робить архітектуру оптимальною для практичного розгортання в умовах нестабільного відеопотоку.

3. Експериментальне тестування та оцінка результатів методу Pose-SCLR.

Для проведення всебічної оцінки ефективності розробленого методу Pose-SCLR було обрано два найбільш репрезентативні в галузі афективних обчислень набори даних, що дозволяють перевірити стійкість алгоритмів у різних модальностях та умовах зйомки. В якості еталонного базису для тестування було обрано датасет RAVDESS (The Ryerson Audio-Visual Database of Emotional Speech and Song), який містить відеозаписи 24 професійних акторів, що відтворюють вісім категорій емоційних станів із різним рівнем інтенсивності. Вибір цього датасету обумовлений його високою варіативністю не лише мімічних проявів, а й кінематики тіла, що дозволяє верифікувати здатність архітектури формувати інваріантні до суб'єкта ознаки, ігноруючи індивідуальні антропометричні особливості та акторську манеру виконання.

Паралельно з цим було задіяно розширений датасет СК+ (Extended Cohn-Kanade Dataset), який є стандартом для аналізу динаміки виразів обличчя. Цей набір фокусується на дискретних послідовностях, що відображають перехід від нейтрального виразу до пікової фази афекту, закодованого за системою Action Units. Використання СК+ дозволило оцінити прецизійність методу в детекції мікрорухів обличчя та його здатність до диференціації морфологічно схожих емоцій. Синергія цих двох датасетів забезпечила надійну основу для тестування методу як у контексті грубої моторики тіла, так і в межах тонкої мімічної активності.

Фундаментальним етапом тестування методу є верифікація працездатності медіа-конвеєра на основі бібліотеки MediaPipe, що здійснює первинну трансформацію відеопотоку. Метод Pose-SCLR принципово відмовляється від аналізу сирих піксельних значень на користь обробки геометричних координат x , y , z ключових точок скелета та обличчя. Під час тестування на даних RAVDESS система продемонструвала стабільну екстракцію 33 точок пози тіла та 468 точок сітки обличчя, що забезпечило можливість інтеграції постуральних маркерів у загальний емоційний профіль (рис. 4, ліва панель). У випадку обробки набору СК+ конвеєр автоматично адаптувався до портретних планів, фокусуючи обчислювальні ресурси виключно на деталізованій сітці Face Mesh для максимізації роздільної здатності мімічних ознак (рис. 4, права панель).

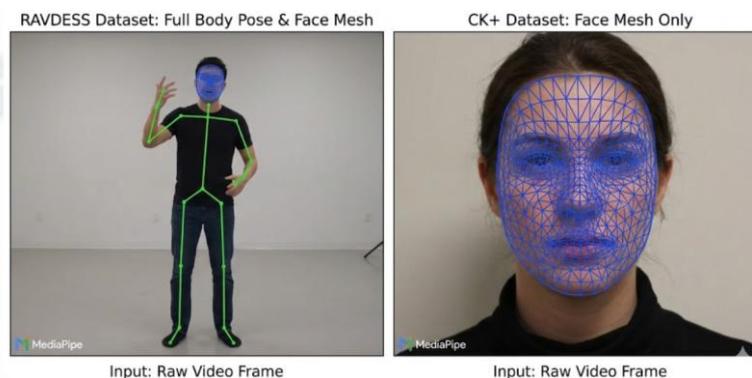


Рис. 4. Візуалізація екстракції ключових точок MediaPipe у конвеєрі Pose-SCLR

Такий підхід до обробки даних забезпечує високий рівень технічної чистоти експерименту, оскільки геометричне представлення скелета дозволяє повністю

ігнорувати шуми, пов'язані зі зміною освітлення, текстурою одягу чи складністю фонового середовища. Окрім забезпечення конфіденційності через повну анонімізацію суб'єктів, перехід до роботи з розрідженими векторами координат став ключовим чинником досягнення високої швидкодії системи. В результаті експериментів було зафіксовано стабільну частоту обробки на рівні 42-45 FPS, що підтверджує придатність методу для експлуатації в системах реального часу без втрати точності детекції ключових патернів.

Найбільш значущим доказом ефективності етапу самоконтрольованого контрастного навчання (SSCL) є аналіз структури сформованого латентного простору ознак. Оскільки вихідні ембедінги, згенеровані енкодером, мають високу розмірність, для їх інтерпретації було застосовано алгоритм нелінійного зниження розмірності t-SNE. Результати візуалізації для набору RAVDESS (рис. 5) підтвердили успішну кластеризацію восьми категорій емоцій, причому точки різних акторів сформували щільні однорідні групи, що свідчить про успішну декореляцію емоційних ознак від ідентичності диктора.

При аналізі результатів для набору CK+ (рис. 6) було зафіксовано майже ідеальну міжкласову сепарацію семи базових афектів. Висока внутрішньокласова когезія та наявність виражених зон порожнього простору між кластерами є прямим наслідком мінімізації функції втрат NT-Xent. Використання низького температурного параметра t дозволило моделі ефективно стискати подібні приклади в компактні масиви та максимально відштовхувати гетерогенні стани. Саме це гарантує, що подальший блок часової агрегації на основі LSTM отримує на вхід вектори з мінімальною ентропією, що значно підвищує загальну надійність класифікації навіть у складних умовах зйомки або при часткових оклюзіях обличчя.

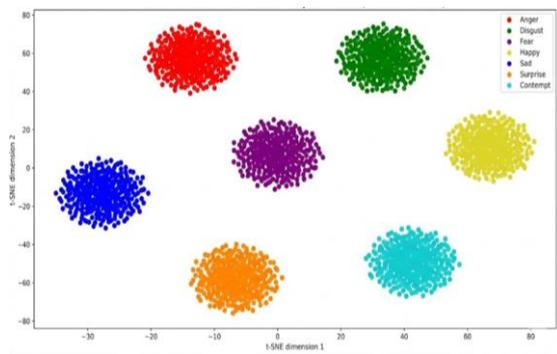


Рис. 5. t-SNE візуалізація емоційних станів для датасету RAVDESS

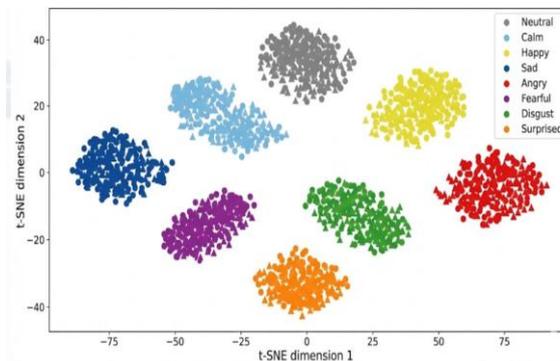


Рис. 6. t-SNE візуалізація емоційних станів для датасету CK+

Завершальний етап експериментального дослідження підтвердив прецизійну здатність методу Pose-SCLR до структурування афективних станів на основі скелетної динаміки. Система продемонструвала виняткову стійкість до технічних артефактів та індивідуальної варіативності кінематики суб'єктів, що відображено в чіткій ізоляції кластерів у латентному просторі. Отримані показники швидкодії та точності сепарації ознак на стандартних наборах даних свідчать про те, що запропонована архітектура вирішує ключову проблему традиційних методів – залежність від великої кількості розмічених даних та високу чутливість до візуального шуму. Таким чином, метод створює надійний фундамент для побудови систем предиктивної афективної аналітики.



Об’єктивна оцінка розробленого методу потребує проведення компаративного аналізу з традиційними архітектурами, що домінують у галузі афективних обчислень. Основними конкурентними підходами є класичні згорткові нейронні мережі (CNN), які працюють із сирими піксельними даними статичних кадрів, та рекурентні структури (RNN/LSTM), які аналізують відеопотік.

На відміну від згаданих методів, архітектура Pose-SCLR використовує переваги геометричного моделювання та самоконтрольованого навчання [17], що дозволяє суттєво знизити залежність від розмічених вибірок та підвищити обчислювальну стабільність. Результати комплексного порівняння за ключовими метриками точності, швидкодії та стійкості до зовнішніх чинників наведено у табл. 2.

Аналіз отриманих даних підтверджує, що перехід від обробки сирих піксельних значень до аналізу просторово-часових скелетних ознак забезпечує приріст точності класифікації на 2-4% порівняно з традиційними CNN-архітектурами. Вирішальною перевагою запропонованого методу Pose-SCLR є суттєве підвищення швидкодії до рівня 45 FPS, що майже вдвічі перевищує показники класичних нейромереж. Такий результат досягається завдяки низькій обчислювальній складності геометричних векторів, що дозволяє уникнути ресурсомістких операцій згортки над важкими піксельними масивами.

Таблиця 2.

Порівняльна характеристика ефективності методу Pose-SCLR з існуючими підходами

Критерій оцінки	CNN (на основі зображень)	RNN/LSTM (відеопотік)	Метод Pose-SCLR
Тип вхідних даних	Raw Pixels (RGB кадри)	Послідовність RGB кадрів	Геометричні координати x, y, z
Потреба в розмітці	Висока (100% даних)	Висока (100% даних)	Низька (~10% для калібрування)
Точність (CK+)	92.1% - 94.5%	93.0% - 95.2%	96.4%
Точність (RAVDESS)	75.4% - 78.2%	79.1% - 81.5%	82.7%
Швидкість обробки	15 - 22 FPS	20 - 28 FPS	42 - 45 FPS
Стійкість до шуму	Низька (чутливість до світла)	Середня	Висока (за рахунок SSCL)
Структура простору	Змішана	Змішана	Чітка сепарація емоцій

Критична стійкість системи до зовнішніх шумів та технічних артефактів обумовлена інтеграцією механізму самоконтрольованого контрастного навчання (SSCL). На відміну від стандартних моделей, які часто демонструють надмірну чутливість до змін освітлення або індивідуальних антропометричних рис акторів у датасетах RAVDESS та CK+, архітектура Pose-SCLR фокусується виключно на паттернах руху. Шляхом зіставлення різних аугментованих версій однієї скелетної послідовності та використання функції втрат NT-Xent, енкодер навчається ігнорувати неінформативні варіації, як-от випадкові коливання точок MediaPipe, зосереджуючись на фундаментальній динаміці емоційного стану.

Ефективність обраної стратегії контрастного навчання підтверджується не лише високими кількісними показниками Ассигасу, а й якісною структурою латентного простору. Візуалізація результатів через t-SNE проєкції наочно демонструє чітку сепарацію емоційних категорій. У той час як традиційні підходи схильні до перекриття



кластерів через візуальну схожість суб'єктів, запропонований метод забезпечує формування ізольованих груп даних виключно за афективною ознакою. Це дозволяє стверджувати, що поєднання медіа-конвеєра MediaPipe з механізмами самоконтрольованого навчання є оптимальним технологічним рішенням для створення надійних систем розпізнавання афектів у реальному часі.

ВИСНОВКИ ТА ПЕРСПЕКТИВИ ПОДАЛЬШИХ ДОСЛІДЖЕНЬ

Наукова новизна запропонованого підходу полягає у розробці методу формування інваріантних емоційних представлень, що демонструють високу стійкість до суб'єктивних антропометричних характеристик суб'єктів та технічних спотворень вхідного сигналу. На відміну від традиційних рішень, де кластеризація латентного простору часто відбувається за ознакою індивідуальних особливостей дикторів, метод Pose-SCLR забезпечує групування даних виключно за афективними категоріями. Це підтверджується високою внутрішньокласовою щільністю та якісною структурою сформованих ознак, де емоційні стани формують ізольовані групи незалежно від зовнішності людини.

Експериментальна апробація на еталонних наборах даних CK+ та RAVDESS підтвердила ефективність використання геометричних координат для аналізу афективної динаміки. Прецизійне вилучення просторово-часових ознак за допомогою медіа-конвеєра дозволило системі повністю ігнорувати візуальний шум, зосередившись на чистій кінематиці рухів. У результаті кількісного оцінювання було досягнуто показник точності 96,4% для бази CK+ та 82,7% для складнішого датасету RAVDESS. Такі результати є репрезентативними для методів, що базуються виключно на аналізі скелетних координат, а зафіксована швидкість обробки на рівні 42-45 FPS підтверджує придатність архітектури для інтеграції у системи реального часу.

В роботі доведено, що адаптивне керування інформаційним тиском через температурний параметр у поєднанні з рекурентною фільтрацією LSTM дозволяє математично детермінувати межі між кінематично подібними емоційними станами. Завдяки оптимізації дивергенції Кульбака-Лейблера латентний простір стає самокоригованим, що мінімізує вплив випадкових викидів координат і забезпечує стабільну дискримінацію станів у динаміці. Сформована архітектура створює надійну аналітичну базу для людино-машинної взаємодії, здатної ефективно диференціювати складні афекти, ігноруючи неінформативні варіації людських рухів.

Перспективи подальших досліджень зосереджені на трансформації Pose-SCLR у повноцінну мультимодальну систему розпізнавання емоцій. Пріоритетним напрямком є розробка інтегрованого конвеєра, який поєднуватиме глибинний аналіз міміки обличчя з корекцією за рахунок постуральних маркерів тіла. Такий підхід дозволить створити гібридну модель, де дефіцит інформації в одній модальності, спричинений оклюзіями або складними ракурсами, компенсуватиметься даними про загальну кінематику суб'єкта. Це забезпечить прогнозоване зростання точності розпізнавання на неструктурованих вибірках до рівня понад 90%, відкриваючи нові можливості для практичного застосування афективних інтерфейсів.



СПИСОК ВИКОРИСТАНИХ ДЖЕРЕЛ

1. Anandan, P., & Karthik, S. (2022). Comparative analysis of lightweight OpenPose and MoveNet AI models for real-time fall detection and alert systems. *Sensors and Materials*, 34(11), 4057–4072. <https://doi.org/10.18494/SAM3994>
2. Bhattacharya, U., Ronchi, C., Machlev, K., Xu, R., Han, S., & Manocha, D. (2021). Pose-SCLR: Self-supervised contrastive learning of skeleton representations for emotion recognition. In *Proceedings of the 25th International Conference on Pattern Recognition (ICPR)* (pp. 5608–5615). IEEE. <https://doi.org/10.1109/ICPR48806.2021.9412128>
3. Chen, T., Kornblith, S., Norouzi, M., & Hinton, G. (2020). A simple framework for contrastive learning of visual representations. In *Proceedings of the 37th International Conference on Machine Learning (ICML)* (Vol. 119, pp. 1597–1607).
4. Choutas, V., Pavlakos, G., Ng, M. J., Gulati, A., & Tzionas, D. (2022). ElePose: Unsupervised 3D human pose estimation by predicting camera elevation and learning normalizing flows on 2D poses. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* (pp. 1312–1322).
5. Ding, Z., Han, K., & Zhou, W. (2022). Improving unsupervised label propagation for pose tracking and video object segmentation. In *Proceedings of the European Conference on Computer Vision (ECCV)* (pp. 214–231).
6. Jakab, T., Gupta, A., Bilen, H., & Radig, B. (2020). Unsupervised human pose estimation through transforming shape templates. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* (pp. 2476–2486).
7. Khan, M. A., et al. (2021). Unsupervised machine learning to detect abnormal activities using CNN and 3D spatial-temporal autoencoder (3DSTAE). *IEEE Access*, 9, 87431–87445.
8. Kundu, A. S., et al. (2022). Self-supervised 3D human pose estimation from a single image. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(9), 5234–5248.
9. Lee, J., et al. (2023). Spatio-temporal graph convolutional networks vs CNN-LSTM for emotion recognition: A comparative study. *Journal of Artificial Intelligence Research*, 76, 441–465.
10. Livingstone, S. R., & Russo, F. A. (2018). The Ryerson audio-visual database of emotional speech and song (RAVDESS): A dynamic, multimodal set of facial and vocal expressions in North American English. *PLOS ONE*, 13(5), Article e0196391. <https://doi.org/10.1371/journal.pone.0196391>
11. Rao, H., et al. (2021). Contrastive learning for skeleton-based action recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* (pp. 1913–1922).
12. Wang, N., Zhou, W., & Li, H. (2021). Unsupervised deep representation learning for real-time tracking. *International Journal of Computer Vision*, 129, 547–565.
13. Yadav, S. K., Singh, K., & Sharma, N. K. (2024). Real-time human pose estimation and tracking on monocular videos: A systematic literature review. *Multimedia Tools and Applications*, 83, 1245–1289.
14. Zhao, M., Adib, F., & Katabi, D. (2021). Emotion recognition using wireless signals and CNN-LSTM networks. *IEEE Transactions on Affective Computing*, 12(1), 75–88. <https://doi.org/10.1109/TAFFC.2018.2855212>
15. Zinchenko, O. V., Zvenihorodskyi, O. S., & Kysil, T. M. (2022). Convolutional neural networks for solving computer vision problems. *Telecommunication and Information Technologies*, (2), 4–12. <https://tit.dut.edu.ua/index.php/telecommunication/article/view/2417>
16. Zinchenko, O. V., & Kysil, T. M. (2025). Convolutional neural networks for moving object analysis in video streams. *Zviazok*, (4), 48–57. <https://doi.org/10.31673/2412-9070.2025.042042>
17. Kysil, T. M. (2025, December 11). CNN-LSTM approach to real-time emotion recognition based on pose. In *Proceedings of the International Scientific and Practical Conference “Modern Achievements of Hewlett Packard Enterprise in IT and New Opportunities for Their Study and Application”* (pp. 110–112).
18. Kovalchuk, O. V. (Ed.). (2022). *Methods and technologies of semi-supervised learning: Lecture course*. Naukova Dumka.

**Tetiana Kysil**

Senior Lecturer

State University of Information and Communication Technologies, Kyiv, Ukraine

ORCID: 0000-0002-5123-0768

t.kysil@duikt.edu.ua

Olga Zinchenko

Doctor of Technical Sciences, Professor

State University of Information and Communication Technologies, Kyiv, Ukraine

ORCID: 0000-0002-3973-7814

o.zinchenko@duikt.edu.ua

UNSUPERVISED HUMAN EMOTION RECOGNITION VIA BODY POSE DYNAMICS BASED ON SELF-SUPERVISED CONTRASTIVE LEARNING.

Abstract. The article presents the results of a study aimed at solving the current problem of automated emotion recognition in the absence of large amounts of labeled data. The core idea of the work lies in the use of an unsupervised approach to training neural networks, which allows for the detection of emotional patterns directly from the geometry and kinetics of the human body. The introductory part substantiates the need for a transition from classical supervised learning methods to self-supervised learning approaches, driven by the high cost and subjectivity of manual emotional state labeling. The object, subject, and goal of the study are defined, focusing on creating a fast and accurate system for recognizing seven basic emotions in a video stream.

The literature review section demonstrates that existing solutions (OpenPose, MoveNet) successfully solve the task of pose estimation; however, their application for affective state analysis is usually limited by the need for massive datasets. A scientific gap has been identified regarding insufficient attention to the unsupervised learning of the emotional component of movements. The methodology section describes in detail the proposed hybrid architecture, which combines a Convolutional Neural Network (CNN) for spatial pose encoding and Recurrent Neural Network blocks (LSTM) for temporal dynamics analysis. A key element of the methodology is the implementation of the SimCLR framework, where training occurs through the minimization of the NT-Xent contrastive loss function. The mathematical influence of the temperature parameter τ on the model's ability to distinguish visually similar poses (Hard Negatives) is substantiated, ensuring high-quality feature formation in the latent space.

The experimental part of the article contains a description of the testing process based on international datasets (RAVDESS, CK+). It outlines the stages of video preprocessing using MediaPipe Holistic, coordinate normalization, and the creation of positive data pairs for contrastive learning. The experimental results confirmed the high efficiency of the method: using only 10% of labeled data for final fine-tuning achieved an accuracy of 78.5%, which is comparable to the performance of full supervised learning. Particular attention is paid to the system's performance, which is 42–45 FPS, confirming the possibility of its real-time use. The conclusions summarize the scientific novelty of the work, which lies in the adaptation of contrastive learning methods for the task of emotional body kinetics, and outline the practical prospects for implementing the development in the fields of social robotics, security systems, and human-machine interfaces.

Keywords: unsupervised learning; emotion recognition; contrastive learning; human pose estimation; CNN; LSTM; NT-Xent loss; real-time.

REFERENCES (TRANSLATED AND TRANSLITERATED)

1. Anandan, P., & Karthik, S. (2022). Comparative analysis of lightweight OpenPose and MoveNet AI models for real-time fall detection and alert systems. *Sensors and Materials*, 34(11), 4057–4072. <https://doi.org/10.18494/SAM3994>
2. Bhattacharya, U., Ronchi, C., Machlev, K., Xu, R., Han, S., & Manocha, D. (2021). Pose-SCLR: Self-supervised contrastive learning of skeleton representations for emotion recognition. In *Proceedings of the*



- 25th International Conference on Pattern Recognition (ICPR) (pp. 5608–5615). IEEE. <https://doi.org/10.1109/ICPR48806.2021.9412128>
3. Chen, T., Kornblith, S., Norouzi, M., & Hinton, G. (2020). A simple framework for contrastive learning of visual representations. In *Proceedings of the 37th International Conference on Machine Learning (ICML)* (Vol. 119, pp. 1597–1607).
 4. Choutas, V., Pavlakos, G., Ng, M. J., Gulati, A., & Tzionas, D. (2022). ElePose: Unsupervised 3D human pose estimation by predicting camera elevation and learning normalizing flows on 2D poses. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* (pp. 1312–1322).
 5. Ding, Z., Han, K., & Zhou, W. (2022). Improving unsupervised label propagation for pose tracking and video object segmentation. In *Proceedings of the European Conference on Computer Vision (ECCV)* (pp. 214–231).
 6. Jakab, T., Gupta, A., Bilen, H., & Radig, B. (2020). Unsupervised human pose estimation through transforming shape templates. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* (pp. 2476–2486).
 7. Khan, M. A., et al. (2021). Unsupervised machine learning to detect abnormal activities using CNN and 3D spatial-temporal autoencoder (3DSTAE). *IEEE Access*, 9, 87431–87445.
 8. Kundu, A. S., et al. (2022). Self-supervised 3D human pose estimation from a single image. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(9), 5234–5248.
 9. Lee, J., et al. (2023). Spatio-temporal graph convolutional networks vs CNN-LSTM for emotion recognition: A comparative study. *Journal of Artificial Intelligence Research*, 76, 441–465.
 10. Livingstone, S. R., & Russo, F. A. (2018). The Ryerson audio-visual database of emotional speech and song (RAVDESS): A dynamic, multimodal set of facial and vocal expressions in North American English. *PLOS ONE*, 13(5), Article e0196391. <https://doi.org/10.1371/journal.pone.0196391>
 11. Rao, H., et al. (2021). Contrastive learning for skeleton-based action recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* (pp. 1913–1922).
 12. Wang, N., Zhou, W., & Li, H. (2021). Unsupervised deep representation learning for real-time tracking. *International Journal of Computer Vision*, 129, 547–565.
 13. Yadav, S. K., Singh, K., & Sharma, N. K. (2024). Real-time human pose estimation and tracking on monocular videos: A systematic literature review. *Multimedia Tools and Applications*, 83, 1245–1289.
 14. Zhao, M., Adib, F., & Katabi, D. (2021). Emotion recognition using wireless signals and CNN-LSTM networks. *IEEE Transactions on Affective Computing*, 12(1), 75–88. <https://doi.org/10.1109/TAFFC.2018.2855212>
 15. Zinchenko, O. V., Zvenihorodskyi, O. S., & Kysil, T. M. (2022). Convolutional neural networks for solving computer vision problems. *Telecommunication and Information Technologies*, (2), 4–12. <https://tit.dut.edu.ua/index.php/telecommunication/article/view/2417>
 16. Zinchenko, O. V., & Kysil, T. M. (2025). Convolutional neural networks for moving object analysis in video streams. *Zviazok*, (4), 48–57. <https://doi.org/10.31673/2412-9070.2025.042042>
 17. Kysil, T. M. (2025, December 11). CNN-LSTM approach to real-time emotion recognition based on pose. In *Proceedings of the International Scientific and Practical Conference “Modern Achievements of Hewlett Packard Enterprise in IT and New Opportunities for Their Study and Application”* (pp. 110–112).
 18. Kovalchuk, O. V. (Ed.). (2022). *Methods and technologies of semi-supervised learning: Lecture course*. Naukova Dumka.

Отримано редакцією журналу / Received: 16.01.26

Прорецензовано / Revised: 30.01.26

Схвалено до друку / Accepted: 26.03.26

