

0.1 Σχετική βιβλιογραφία

Σε αυτή την ενότητα, παρουσιάζουμε ορισμένες από τις πιο συναφείς μεθόδους για την εργασία μας και άλλες που μελετήθηκαν για τον σχεδιασμό αυτής της προσέγγισης. Οι μέθοδοι αυτές χωρίζονται σε δύο ενότητες *Αναγνώριση Δραστηριότητας* και *Εντοπισμός Δραστηριότητας*. Το πρώτο μέρος αναφέρεται σε κλασικές μεθόδους ταξινόμησης δράσης που εισήχθησαν μέχρι πρόσφατα και το δεύτερο μέρος, αντίστοιχα, σε πρόσφατες μεθόδους εντοπισμού της δράσης.

0.1.1 Αναγνώριση Δραστηριότητας

Οι πρώτες προσεγγίσεις για την κατάταξη της δράσης αποτελούνταν από δύο βήματα α) αρχικά υπολογισμός σύνθετων ‘χειροποίητων’ χαρακτηριστικών από ακατέργαστα καρέ βίντεο και β) εκπαίδευση ενός ταξινομητή με βάση αυτά τα χαρακτηριστικά. Αυτά τα χαρακτηριστικά μπορούν να διαχωριστούν σε 3 κατηγορίες: 1) προσεγγίσεις χωροχρονικού όγκου (space-time volume), 2) τροχιές (trajectories) και 3) χωροχρονικά χαρακτηριστικά. Για τις μεθόδους χωροχρονικού όγκου, η προσέγγιση είναι η εξής: Με βάση τα training βίντεο, το σύστημα συνάπτει ένα μοντέλο τρισδιάστατου χωροχρόνου, συνενώνοντας δυσδιάστατες εικόνες (διάσταση $x-y$) κατά τη διάρκεια του χρόνου (διάσταση t ή z), για την αναπαράσταση κάθε δράσης. Όταν το σύστημα δέχεται ένα βίντεο που δεν έχει ετικέτα, κατασκευάζει μια τρισδιάστατη χωροχρονική αναπαράσταση που αντιστοιχεί σε αυτό το βίντεο. Αυτό η νέα τρισδιάστατη αναπαράσταση, στη συνέχεια, συγκρίνεται με κάθε μοντέλο 3D χωροχρόνου, συγκρίνοντας την ομοιότητα στο σχήμα και την εμφάνιση μεταξύ αυτών των δύο χωροχρονικών όγκων. Το σύστημα εξάγει την κατηγορία του άγνωστου βίντεο, αντιστοιχώντας την με αυτήν της δράσης με την υψηλότερη ομοιότητα. Επιπλέον, υπάρχουν διάφορες παραλλαγές των χωροχρονικών αναπαράστασεων. Αντί της αναπαράστασης space-time volume, το σύστημα μπορεί να αναπαριστά τη κάθε δράση ως τροχιές σε χωροχρονικές διαστάσεις ή ακόμη περισσότερο, η ενέργεια μπορεί να αναπαρασταθεί ως ένα σύνολο χαρακτηριστικών που εξάγονται από τον χωροχρονικό όγκο ή τις τροχιές. Οι ‘καθαρές’ χωροχρονικές αναπαραστάσεις περιλαμβάνουν μεθόδους σύγκρισης των περιοχών προσκήνιου ενός ατόμου (δηλ. σιλουέτες) όπως Bobick and Davis 2001, συγκρίνοντας όγκους σε σχέση με επιφάνεια τους όπως οι Shechtman and Irani 2005. Η μέθοδος Ke, Sukthankar, and Hebert 2007 χρησιμοποιεί oversegmented όγκους, αυτομάτως υπολογίζοντας ένα σύνολο τμημάτων τρισδιάστατου όγκου XYT που αντιστοιχεί σε έναν κινούμενο άνθρωπο. Rodriguez, Ahmed, and Shah 2008 πρότειναν φίλτρα για να αποτυπώνουν τα χαρακτηριστικά του χωροχρονικού όγκου, προκειμένου να τα ταιριάζουν πιο αξιόπιστα και αποδοτικά. Από την άλλη πλευρά, οι προσεγγίσεις με βάση την τροχιά περιλαμβάνουν την αναπαράσταση μιας ενέργειας ως σύνολο 13 κοινών διαδρομών (Sheikh, Sheikh, and Shah 2005) ή τη χρήση ενός συνόλου XYZT-διαστάσεων κοινών τροχιών που λαμβάνονται από κινούμενες κάμερες (Yilmaz and Shah 2005). Τέλος, διάφορες μέθοδοι χρησιμοποιούν τοπικά χαρακτηριστικά που εξάγονται από χωροχρονικούς όγκους τριών διαστάσεων, όπως η εξαγωγή τοπικών χαρακτηριστικών σε κάθε καρέ του βίντεο και η ένωση του χρονικά (Chomat and Crowley 1999; Zelnik-Manor and Irani 2001; Blank

et al. 2005, η εξαγωγή αραιών χωροχρονικών τοπικών σημείων ενδιαφέροντος από τρισδιάστατους όγκους (Laptev and Lindeberg 2003; Dollar et al. 2005; Niebles, Wang, and Li 2006; Alper Yilmaz and Mubarak Shah 2005; Ryoo and Aggarwal 2006) Οι προσεγγίσεις αυτές κατέστησαν την επιλογή των χαρακτηριστικών σημαντικό παράγοντα για την απόδοση του δικτύου. Αυτό συμβαίνει επειδή οι διαφορετικές κατηγορίες δράσεων μπορεί να διαφέρουν δραματικά από την άποψη της εμφάνισής τους και των μοτίβων κίνησης. Ένα άλλο πρόβλημα ήταν ότι οι περισσότερες από αυτές τις προσεγγίσεις κάνουν υποθέσεις, υπό τις οποίες το βίντεο λήφθηκε λόγω προβλημάτων όπως το γεμάτο φόντο, γωνιές κάμερας κλπ. Μια ανασκόπηση των τεχνικών, που χρησιμοποιούνταν μέχρι το 2011, παρουσιάζεται στο Aggarwal and Ryoo 2011.

Τα πρόσφατα αποτελέσματα σε βαθιές αρχιτεκτονικές και ειδικά στον τομέα της ταξινόμησης εικόνας έδωσε κίνητρο στους ερευνητές να εκπαιδεύσουν δίκτυα CNN για το πρόβλημα της αναγνώρισης δράσης. Η πρώτη σημαντική απόπειρα έγινε από τους Karpathy et al. 2014. Σχεδίασαν την αρχιτεκτονική τους με βάση το καλύτερο CNN στον διαγωνισμό ImageNet. Εξερευνούν διάφορες μεθόδους για τη σύντηξη των χωροχρονικών λειτουργιών χρησιμοποιώντας δυσδιάστατες διαδικασίες κυρίως και τρισδιάστατη συνέλιξη μόνο σε αργή σύντηξη. Οι Simonyan and Zisserman 2014 χρησιμοποίησαν 2 "NNs, ένα για χωρικές πληροφορίες και ένα για οπτική ροή και τα συνδύασαν με τη χρήση της καθυστερημένης σύντηξης. Δείχνουν ότι η εξόρυξη χωρικού περιεχομένου από τα βίντεο και περιεχόμενο κίνησης από την οπτική ροή μπορεί να βελτιώσει σημαντικά την ακρίβεια της αναγνώρισης της δράσης. Οι Feichtenhofer, Pinz, and Zisserman 2016 επέκτειναν αυτή την προσέγγιση με τη χρήση πρώιμης σύντηξης στο τέλος των convolutional layers αντί για καθυστερημένης σύντηξης, η οποία λαμβάνει χώρα στο τελευταίο επίπεδο του δικτύου. Πάνω σ' αυτό, χρησιμοποίησαν ένα δεύτερο δίκτυο για το χρονικό περιεχόμενο το οποίο συνδέουν με το άλλο δίκτυο με χρήση της καθυστερημένης σύντηξης. Επιπλέον, οι Wang et al. 2016 στήριξαν την μέθοδο τους σε αυτήν που πρότειναν οι Simonyan and Zisserman 2014. Ασχολούνται με το πρόβλημα του την εύρεσης χρονικού περιεχομένου και εκπαιδεύουν το δίκτυο τους, παρέχοντας του λίγα δείγματα. Η προσέγγισή τους, την οποία ονόμασαν Temporal Segment Network (TSN), διαχωρίζει το βίντεο εισόδου σε K τμήματα και ένα σύντομο απόσπασμα από κάθε τμήμα επιλέγεται για ανάλυση. Στη συνέχεια, συνδέουν το εξαγόμενο χωροχρονικό περιεχόμενο, πραγματοποιώντας τελικά την πρόβλεψή τους. Πιο πρόσφατα, οι Zhang et al. 2016 και οι Zhu et al. 2017 χρησιμοποίησαν την two-stream , επίσης. Οι Zhang et al. 2016 αντικατέστησαν την οπτική ροή με ένα διάνυμα κίνησης που μπορεί να ληφθεί απευθείας από τα συμπίεσμένα βίντεο χωρίς επιπλέον υπολογισμό και το τροφοδοτούν στο δίκτυο. Οι Zhu et al. 2017 εκπαιδεύσαν ένα CNN για τον υπολογισμό της οπτικής ροής, καλώντας το MotionNet και χρησιμοποίησαν ένα CNN ως χρονικό stream για προβάλλουν τις πληροφορίες κίνησης έργου σε κατηγορίες δράσεων. Τέλος χρησιμοποιούν Την καθυστερημένη σύντηξη μέσω της μέσης τιμής με βάση των σκορ πρόβλεψης των χρονικών και χωρικών stream. Από την άλλη πλευρά, μια νέα προσέγγιση εισήχθη από τους Γκρδθαφ και Ραμαναν 2017 ενσωματώνοντας χάρτες προσοχής με σκοπό να βελτιώσουν σημαντικά την απόδοση της αναγνώρισης δράσης.

Ορισμένες άλλες μέθοδοι περιλάμβαναν ένα δίκτυο RNN ή LSTM για την ταξινόμηση κάδων οι Donahue et al. 2017, οι Joe Yue-Hei Ng et al. 2015 και οι Ma et al. 2017. Οι Donahue et al. 2017 αντιμετωπίζουν τη πρόκληση των μεταβλητών μήκη των ακολουθιών εισόδου και εξόδου, εκμεταλλευόμενοι τα convolutional layers και μεγάλου εύρους χρονικές αναδρομές (recursions). Προτείνουν ένα Long-term Recurrent Convolutional Network (LRCN), το οποίο είναι ικανό να αντιμετωπίσει τις εργασίες αναγνώρισης, λεζάντας εικόνας και περιγραφής βίντεο. Για να ταξινομήσουν μια δεδομένη ακολουθία καρέ, το LRCN λαμβάνει αρχικά ως είσοδο ένα καρέ, και πιο συγκεκριμένα τα κανάλια RGB και την οπτική ροή του, και προβλέπει μια ετικέτα. Μετά από αυτό, εξάγει την κλάση του βίντεο μέσω του μέσου όρου των πιθανοτήτων των ετικετών, επιλέγοντας την πιο πιθανή κλάση. Οι Joe Yue-Hei Ng et al. 2015 πρώτα διερευνούν διάφορες προσεγγίσεις για χρονική ομαδοποίησης (temporal pooling) των χαρακτηριστικών. Αυτές οι τεχνικές περιλαμβάνουν τον χειρισμό καρέ βίντεο ξεχωριστά από 2 αρχιτεκτονικές CNN: είτε απ' το AlexNet είτε απ' το GoogleNet, και αποτελούνται από πρώιμη σύντηξη, καθυστερημένη σύντηξη και ενός συνδυασμού αυτών. Επιπλέον, προτείνουν ένα RNN προκειμένου να εξετάσουν τα βίντεο κλιπ ως ακολουθίες ενεργοποιήσεων NN. Το προτεινόμενο LSTM λαμβάνει ως είσοδο την έξοδο του τελικού CNN layer για κάθε συνεχόμενο καρέ και μετά από 5 LSTM layers και χρησιμοποιώντας έναν softmax ταξινομητή, προτείνει μία ετικέτα. Για την ταξινόμηση του βίντεο, επιστρέφουν μια ετικέτα μετά το τελευταίο βήμα, εφαρμόζουν max-pooling στις προβλέψεις στην διάσταση του χρόνου, αθροίζουν τις προβλέψεις στην διάσταση του χρόνου και επιστρέφουν το μέγιστο ή έναν γραμμικό συνδυασμό με βάρη των προβλέψεων υπό έναν παράγοντα g, τα αθροίζουν και επιστρέφουν το μέγιστο. Έδειξαν ότι όλες οι προσεγγίσεις είναι 1% διαφορετικές με προκατάληψη για τη χρήση των προβλέψεων με βάρη για την υποστήριξη της ιδέας ότι το LSTM γίνεται προοδευτικά πιο ενημερωμένο. Τελευταίοι αλλά όχι λιγότερο σημαντικοί, οι Ma et al. 2017 χρησιμοποίησαν ένα two-stream ConvNet για εξαγωγή χαρακτηριστικών και είτε ένα LSTM ή convolutional πάνω από το χρονικώς κατ'ασχευασμένους πίνακες χαρακτηριστικών για τη σύντηξη χωρικών και χρονικών πληροφοριών. Χρησιμοποιούν ένα ResNet-101 για την εξόρυξη χαρτών ενεργοποίησης τόσο για χωρικές όσο και για χρονικές ροές. Χωρίζουν το βίντεο σε διάφορα τμήματα, όπως έκαναν οι Wang et al. 2016, και χρησιμοποίησαν ένα επίπεδο temporal pooling για την εξαγωγή διακεκριμένων χαρακτηριστικών. Αφού λάβουν αυτά τα χαρακτηριστικά, το LSTM εξάγει ενσωματωμένες δυνατότητες από όλα τα τμήματα.

Επιπλέον, οι Tran et al. 2015 διευρένησαν τα 3D Convolutional δίκτυα (Ji et al. 2013) και εισήγαγαν το C3D δίκτυο που έχει 3D convolutional layers με πυρήνες $3 \times 3 \times 3$. Αυτό το δίκτυο είναι σε θέση να μοντελοποιήσει την εμφάνιση και την κίνηση ταυτόχρονα χρησιμοποιώντας τρισδιάσττες συνελίξεις και μπορεί να χρησιμοποιηθεί ως εξαγωγέας χαρακτηριστικών. Συνδυάζοντας την αρχιτεκτονική δύο ροών και τις τρισδιάσττες συνελίξεις οι Άρρεϊρα και Ζισσερμαν 2017 πρότειναν το δίκτυο I3D. Πάνω σ' αυτό, οι δημιουργοί τονίζουν τα πλεονεκτήματα της μεταφοράς μάθησης για την εργασία της αναγνώρισης επαναλαμβάνοντας τα δυοδιάστατα προ-εκπαιδευμένα βάρη στην 3η διάσταση. Οι Hara, Kataoka, and Satoh 2017 πρότειναν ένα δίκτυο 3D ResNet για την αναγνώριση δράσης με βάση

τα Residual δίκτυα (ResNet)(He et al. 2016) και διερευνούν την απόδοση των δικτύων ResNet με 3D Convolutional πυρήνες. Από την άλλη, οι Diba et al. 2017 βάσισαν την προσέγγισή τους στα DenseNets (Huang et al. 2017) και επέκτειναν την αρχιτεκτονική του DenseNet χρησιμοποιώντας τρισδιάστατα φίλτρα και pooling πυρήνες αντί για δισδιάστατους, ονομάζοντας αυτή την προσέγγιση ως DenseNet3D. Επιπλέον, εισάγουν το Layer χρονικής μετάβασης (TTL), το οποίο συνενώνει χρονικά χάρτες χαρακτηριστικών που εξάγονται σε διαφορετικά χρονικά βήθη και αντικαθιστά το επίπεδο μετάβασης του DenseNet. Παράλληλα, οι Diba et al. 2018 εισήγαγαν ένα νέο χρονικό layer το οποίο μοντελοποιεί μεταβλητούς χρονικούς πυρήνες συνέλιξης. Τελευταίοι αλλά εξίσου σημαντικοί, οι Tran et al. 2018 πειραματίστηκαν με διάφορες υπόλοιπες αρχιτεκτονικές Residual δικτύου χρησιμοποιώντας συνδυασμούς 2D και 3D ροζολυτιοναλ Λαψερ. Σκοπός τους είναι να δείξουν ότι η 2D χωρική συνέλιξη ακολουθούμενη από 1D χρονική συνέλιξη επιτυγχάνει state-of-the-art αποτελέσματα, ονομάζοντας αυτού του τύπου το layer ως $R(2 + 1)D$. Πρόσφατα οι Guo et al. 2018 πρότειναν ένα framework που μπορεί να μάθει να αναγνωρίζει μια προηγούμενης αθέατη 3D κλάση δράσης με λίγα μόνο παραδείγματα εκμεταλλευόμενο την εγγενή δομή των 3D δεδομένων μέσω μιας γραφικής αναπαράστασης. Ακόμα πιο λεπτομερή παρουσίαση των τεχνικών αναγνώρισης δράσης που χρησιμοποιήθηκαν μέχρι το 2018 πραγματοποιήθηκε από τους Kong and Fu 2018.

0.1.2 Εντοπισμός Δραστηριότητας

Όπως προαναφέρθηκε, ο εντοπισμός δράσης μπορεί να θεωρηθεί ως προέκταση του προβλήματος εντοπισμού αντικειμένων. Αντί να εξάγουμε δισδιάστατα πλαίσια οριοθέτησης σε μία μόνο εικόνα, ο στόχος των συστημάτων εντοπισμού δράσης είναι να εξάγουν action tubes, τα οποία είναι ακολουθίες πλαισίων οριοθέτησης που περιέχουν μια ενέργεια που εκτελέστηκε. Έτσι, υπάρχουν διάφορες προσεγγίσεις, συμπεριλαμβανομένου συνήθως ενός δικτύου ανιχνευτή αντικειμένων και ενός ταξινόμητη.

Οι πρώτες προσεγγίσεις ανιχνευσης αντικειμένων περιλάμβαναν την επέκταση ενός αλγορίθμου πρότασης αντικειμένων σε 3 διαστάσεις. Οι Tian, Sukthankar, and Shah 2013 επέκτειναν τα παραμορφώσιμα (deformable) μοντέλα (Felzenszwalb et al. 2010) με το να αντιμετωπίζουν τις δράσεις ως χωροχρονικά μοτίβα και δημιούργησαν ένα παραμορφώσιμου μοντέλο για κάθε δράση. Οι Jain et al. 2014 εισήγαγαν την έννοια των tubelets, γνωστά και ως ακολουθίες πλαισίων οριοθέτησης και βάσισαν τη μέθοδό τους σε επιλεκτικό αλγόριθμο αναζήτησης (Uijlings et al. 2013), επεκτείνοντας τα superpixels σε supervoxels για την παραγωγή χωροχρονικών σχημάτων. Από την άλλη, οι Oneata et al. 2014 επέκτειναν μια τυχαιοποιημένη διαδικασία συγχώνευσης superpixels που χρησιμοποιούταν για που χρησιμοποιούνταν για προτάσεις αντικειμένων, όπως παρουσιάστηκαν από τους Manen, Guillaumin, and Gool 2013. Οι Yu and Yuan 2015 πρώτα προτείνουν πλαίσια οριοθέτησης για κάθε καρέ με χρήση ενός ανιχνευτή ανθρώπου και κίνησης, ενώ, στη συνέχεια, με τη επιλογή των καλύτερων σε σκορ κουτιών, πρότειναν έναν άπληστο συνδυαστικό αλγόριθμο με τη διατύπωση την εργασίας σύνδεσης ως πρόβλημα μέγιστης κάλυψης. Οι Gemert et al. 2015 παράγουν χωροχρονικές προ-

τάσεις κατευθείαν από τις πυκνές τροχειές, οι οποίες επίσης χρησιμοποιήθηκαν για ταξινόμηση. Οι Chen and Corso 2015 δημιουργούν ένα γράφημα χωροχρονικής τροχιάς και επιλέγουν προτάσεις δράσεων που βασίζονται μόνο στην εσκεμμένη κίνηση που εξάγεται από το γράφημα. Οι Soomro, Idrees, and Shah 2015 διαχωρίζουν τα τμήματα βίντεο σε supervoxels και χρησιμοποιούν το περιεχόμενο τους ως χωρική σχέση μεταξύ των supervoxels σε σχέση με την δράση του προσκηνίου. Δημιουργούν ένα γράφημα για κάθε βίντεο, όπου τα υπεροξέλες σχηματίζουν τους κόμβους και οι κατευθυνμένες άκρες απεικονίζουν τις χωρικές σχέσεις μεταξύ τους. Κατά τη διάρκεια των δοκιμών, κάνουν μια βόλτα στο περιβάλλον, όπου κάθε βήμα καθοδηγείται από τις σχέσεις περιβάλλοντος κατά τη διάρκεια της εκπαίδευσης, με αποτέλεσμα μια κατανομή πιθανότητας μιας δράσης για όλα τα υπεροξέλες. Οι Mettes, Gemert, and Snoek 2016 αντί για τοποθέτηση πλασιών σε όλα τα καρέ των βίντεο, σχολίασαν σημεία σε ένα αραιό υποσύνολο καρέ του βίντεο και χρησιμοποίησαν προτάσεις που λαμβάνονται μέσω ενός μέτρου επικάλυψης μεταξύ των προτάσεων δράσης και των σημείων. Οι Behl et al. 2017 ασχολούνται με την ανίχνευση και τον εντοπισμό ενεργειών σε πραγματικό χρόνο μέσω της λήψης προτάσεων δράσης ανά καρέ και την πρόταση ενός αλγορίθμου σύνδεσης που είναι σε θέση να κατασκευάσει και να ενημερώνει τα action tubes ανά καρέ. Πιο πρόσφατα, οι Soomro and Shah 2017 προσπάθησαν να ασχοληθούν με το πρόβλημα της ανίχνευσης και τον εντοπισμό δράσης χωρίς επίβλεψη. Η προσέγγισή τους περιελάμβανε αρχικά την εξόρυξη κατακερματισμένων supervoxel και στη συνέχεια την ανάθεση ένα βάρους σε κάθε supervoxel. Με την εξαγωγή supervoxels, δημιουργούν ένα γράφημα και στη συνέχεια χρησιμοποιούν μια διακριτική clustering προσέγγιση εκπαιδεύεται ένας ταξινομητής.

Η εισαγωγή του R-CNN (Girshick et al. 2014) κατάφερε σημαντικές βελτιώσεις στην απόδοση των δικτύων εντοπισμού αντικειμένων. Αυτή η αρχιτεκτονική, πρώτον, προτείνει περιοχές στην εικόνα που είναι πιθανό να περιέχουν κάποιο αντικείμενο και στη συνέχεια, τα ταξινομεί χρησιμοποιώντας ένα SVM. Εμπνευσμένοι από αυτή την αρχιτεκτονική, οι Gkioxari and Malik 2015 σχεδιάσαν ένα δίκτυο RCNN 2-stream για να προτείνει προτάσεις δράσεων για κάθε καρέ, ένα stream για το επίπεδο καρέ και ένα για την οπτική ροή. Στη συνέχεια, τα συνδέουν χρησιμοποιώντας τον αλγόριθμο σύνδεσης Viterbi. Οι Weinzaepfel, Harchaoui, and Schmid 2015 επεκτείνουν αυτή την προσέγγιση, εκτελώντας προτάσεις στο επίπεδο καρέ και χρησιμοποιώντας ένα τραςκερ για τη σύνδεση των προτάσεων αυτών χαρακτηριστικά της χωρικής και οπτικής ροής. Επίσης, η μέθοδός τους εκτελεί χρονικό εντοπισμό μέσω της χρήσης ενός συρόμενου παράθυρου πάνω από τα εντοπισμένα tubes.

Η εισαγωγή του Faster RCNN (Ren et al. 2017) συνήσφερε πολύ τη βελτίωση της απόδοσης των δικτύων εντοπισμού δράσης. Οι Peng and Schmid 2016 και Saha et al. 2016 χρησιμοποιούν το Faster R-CNN αντί για το RCNN για προτάσεις σε επίπεδο καρέ, χρησιμοποιώντας το RPN για εικόνες RGB και οπτικής ροής. Αφού λάβουν χωρικές προτάσεις και προτάσεις κίνησης, οι Peng and Schmid 2016 τις συγχωνεύουν και από κάθε προτεινόμενη ROI, παράγουν 4 ROIs για να επικεντρωθούν σε συγκεκριμένο μέρος του σώματος του δρώντα. Μετά από αυτό, συνδέουν την πρόταση χρησιμοποιώντας τον αλγόριθμο Viterbi για κάθε κλάση και εκτελούν χρονικό εντοπισμό χρησιμοποιώντας ένα συρόμενο παράθυρο,

με πολλαπλές χρονικές κλίμακες και διασκελισμό κάνοντας χρήση μιας μεθόδου μέγιστης υποσυστοιχίας. Από την άλλη, οι Saha et al. 2016 εκτελούν, επίσης, ταξινόμηση σε επίπεδο καρτέ. Μετά από αυτό, η μέθοδός τους εκτελεί σύντηξη με βάση έναν συνδυασμό της εμφάνισης και της κίνησης με βάση τις προτάσεις και την βαθμολογία αλληλεπικάλυψης. Τέλος, η χρονική προσαρμογή λαμβάνει χώρα χρησιμοποιώντας δυναμικό προγραμματισμό. Παράλληλα, οι Weinzaepfel, Martin, and Schmid 2016 χρησιμοποιούν το Faster RCNN για την εξαγωγή ανθρώπινων tubes από βίντεο που εστιάζουν στο πρόβλημα του ασθενώς εποπτευόμενου εντοπισμού δράσης. Στη συνέχεια, χρησιμοποιώντας πυκνές τροχιές και μια multi-fold Multiple Instance Learning προσέγγιση (Cinbis, Verbeek, and Schmid 2016) εκπαιδεύουν ένα ταξινομητή. Οι Mettes and Snoek 2017 εισήγαγαν μια μέθοδο για zero-shot Εντοπισμού δράσης. Η προσέγγισή τους περιλαμβάνει την βαθμολόγηση των προτεινόμενων action tubes σύμφωνα με τις αλληλεπιδράσεις μεταξύ των ατόμων που δρουν και αντικειμένων. Χρησιμοποίησαν το Faster-RCNN, στο πρώτο βήμα, για την ανίχνευση τόσο των ανθρώπων που δρουν όσο και των αντικειμένων και μετά, χρησιμοποιώντας χωρικές σχέσεις μεταξύ τους, συνδέουν τα προτεινόμενα πλαίσια στον άξονα του χρόνου βασιζόμενοι στην zero-shot πιθανότητα της παρουσίας των ατόμων, συναφών αντικειμένων γύρω από αυτούς και τις αναμενόμενες χωρικές σχέσεις μεταξύ αντικειμένων και ανθρώπων που δρουν. Επιπλέον οι He et al. 2018 πρότειναν το Tube Proposal Network (TPN) για τη δημιουργία ανεξαρτήτου κλάσης προτάσεις tubelet, οι οποίες χρησιμοποιούν το Faster R-CNN για να λάβουν δισδιάστατες προτάσεις περιοχών και έναν αλγόριθμο σύνδεση για τη σύνδεση των tubelets με τις προτάσεις των περιοχών. Πιο πρόσφατα, οι Girdhar et al. 2018 πρότειναν μια μέθοδο για εντοπισμό δράσεων στο σύνολο δεδομένων AVA (Gu et al. 2018) συνδυάζοντας τις αρχιτεκτονικές των I3D (Carreira and Zisserman 2017) και Faster RCNN. Χρησιμοποιούν μπλοκ του I3D για την λήψη αναπαράστασης βίντεο και το RPN του Faster-RCNN για να προτείνει προτάσεις «ανθρώπου» για το κεντρικό πλαίσιο.

Παράλληλα με αυτά, οι **singh2016online** και **kalogeiton17iccv:hal-01519812** σχεδίασαν τα δίκτυα τους με βάση το Single Shot Multibox Detector **DBLP:journals/corr/LiuAESR15**). Οι **singh2016online** δημιούργησαν ένα χωροχρονικό δίκτυο πραγματικού χρόνου. Για να λειτουργεί το δίκτυο τους σε πραγματικού χρόνου, **singh2016online** πρότειναν έναν νέο και αποδοτικό αλγόριθμο με την προσθήκη πλαισίων σε tubes σε κάθε καρτέ, εάν επικαλύπτονται περισσότερο από ένα κατώφλιο, ή εναλλακτικά, τερματίζουν το action tube εάν για k καρτέ δεν προσθέθηκε κανένα πλαίσιο. Οι **kalogeiton17iccv:hal-01519812** σχεδίασαν ένα δίκτυο δύο ροών, το οποίο κάλεσαν ACT-detector, και εισήγαγαν τα κυβικά (cuboids) anchors. Για K καρτέ, και για τα δύο δίκτυα, οι **kalogeiton17iccv:hal-01519812** εξάγουν χωρικά χαρακτηριστικά σε επίπεδο καρτέ, στη συνέχεια, τα στοιβάζουν. Τέλος, με τη χρήση των κυβικών anchors, το δίκτυο εξάγει tubelets, με τις αντίστοιχες βαθμολογίες κατάταξης και στόχους παλινδρόμησης. Για τη σύνδεση των tubelets, οι **kalogeiton17iccv:hal-01519812** ακολουθούν τα ίδια βήματα με τους **singh2016online**. Για χρονική εντοπισμό, χρησιμοποιούν μία προσέγγιση χρονικής εξομάλυνσης.

Πιο πρόσφατα, το δίκτυο YOLO (**DBLP:journals/corr/RedmonDGF15**) έγινε η έμπνευση για τους **DBLP:journals/corr/abs-1903-00304** και τους

DBLP:journals/corr/abs-1802-08362. Στην προσέγγιση που προτάθηκε από τους **DBLP:journals/corr/abs-1903-00304**, οι έννοιες της εξέλιξης και τού ποσοστού προόδου εισήχθησαν. Εκτός από την πρόταση πλαισίων οριοθέτησης σε επίπεδο καρέ, χρησιμοποιούν το YOLO μαζί με έναν ταξινομητή RNN για να εξάγουν χρονικές πληροφορίες για τις προτάσεις. Με βάση αυτές τις πληροφορίες, δημιουργούν action tubes, χωρίζοντας τα σε κλάσεις. Ορισμένες άλλες προσεγγίσεις περιλαμβάνουν εκτίμηση πόζας αυτή των **DBLP:journals/corr/abs-1802-09232**.

Πρότειναν μια μέθοδο υπολογισμού των δισδιάστατων και τρισδιάστατων πόζων και στη συνέχεια εκτέλεσαν ταξινόμηση δράσεων. Χρησιμοποιούν το διαφορίσιμο Soft-argmax για την εκτίμηση των 2D και 3D αρθρώσεων, επειδή η συνάρτηση argmax δεν είναι διαφορίσιμη. Στη συνέχεια, για T παρακείμενες δημιουργούν μια απεικόνιση εικόνας με το χρόνο και τις N_j αρθρώσεις ως $x - y$ άξονες, έχοντας 2 κανάλια για την 2D πόζα και 3 για την 3D πόζα. Χρησιμοποιούν Convolutional Layers για να παράγουν χάρτες θερμότητας δράσης και στη συνέχεια χρησιμοποιώντας max plus min pooling και την συνάρτηση softmax εκτελούν ταξινόμηση δράσης. Οι **DBLP:journals/corr/ZolfaghariOSB17** πρότειναν μια αρχιτεκτονική τριών ροών που περιλαμβάνει 2D πόζα, οπτική ροή και πληροφορίες RGB. Αυτά τα streams ενώνονται μέσω του μοντέλου της αλυσίδας Μάρκοφ. Επιπλέον, οι **8237881** πρότειναν μια αρχιτεκτονική με τη χρήση ενός χρονικού convolutional δικτύου παλινδρόμησης, για να πιάνουν την μακροπρόθεσμη εξάρτηση και πληροφορίες μεταξύ γειτονικών καρέ και ένα χωρικό δίκτυο παλινδρόμησης, για προτάσεις ανά καρέ. Χρησιμοποιούν μεθόδους παρακολούθησης και δυναμικού προγραμματισμού για τη δημιουργία προτάσεων δράσης.

Τα περισσότερα από τα προαναφερθέντα δίκτυα χρησιμοποιούν ανά καρέ χωρικές προτάσεις και εξάγουν τις χρονικές τους πληροφορίες υπολογίζοντας την οπτική ροή. Από την άλλη οι **DBLP:journals/corr/SahaSC17** σχεδίασαν μια αρχιτεκτονική η οποία περιλαμβάνει προτάσεις σε επίπεδο τμήματος βίντεο, το οποίο σημαίνει περισσότερα από ένα καρέ ταυτόχρονα. Οι **DBLP:journals/corr/SahaSC17** πρότειναν μια 3D-RPN αρχιτεκτονική που είναι σε θέση να δημιουργήσει και να ταξινομήσει τρισδιάσττες προτάσεις αποτελούμενες από 2 συνεχόμενα καρέ. Επίσης, πρότειναν έναν αλγόριθμο σύνδεσης, τροποποιώντας αυτόν που πρότειναν οι Saha et al. 2016. Πάνω σ' αυτό, οι **DBLP:journals/corr/HouCS17** σχεδίασαν μια αρχιτεκτονική για τη δημιουργία προτάσεων δράσης για περισσότερα από 2 καρέ, καλώντας το μοντέλο τους Tube CNN (T-CNN). Στην προσέγγισή τους, το επίπεδο του τμήματος βίντεο σημαίνει ότι ολόκληρο το βίντεο χωρίζεται κλιπ βίντεο ίδιου αριθμού καρέ και με τη χρήση του C3D για την εξόρυξη χαρακτηριστικών, επιστρέφουν χωροχρονικές προτάσεις. Μετά την λήψη των προτάσεων, οι **ΔΒΛΠ:θουρναλς/ςορρ/Ηου"Σ17** συνδέουν τις tube προτάσεις τους με έναν αλγόριθμο στηριζόμενος στην πιθανότητα ύπαρξης δράσης και την επικάλυψη μεταξύ των tubes. Τέλος, η λειτουργία ταξινόμησης λαμβάνει χώρα για τα συνδεδεμένα action tubes.

Βιβλιογραφία

- [1] Ο. Ήοματ και Θ. Α. Ήρωλεψ. ‘Προβαβιλιστις ρεσογνιτιον οφ αστιτψ υσιγγ λοσαλ αππεαρανσε’. Στο: *Προσεεδινγς. 1999 IEEE δμπτυερ Σοσιετψ δνφερενς ον δμπτυερ ίσιον ανδ Παττερν Ρεσογνιτιον* (ατ. Νο ΠΡ00149). Τόμ. 2. 1999, 104–109 έλ. 2. ΔΟΙ: 10.1109/΄΄ΠΡ.1999.784616.
- [2] Α. Φ. Βοβιςκ και Θ. Ω. Δαις. ‘Τηε ρεσογνιτιον οφ ημμαν μοεμεντ υσιγγ τεμποραλ τεμπλατες’. Στο: *IEEE Τρανσαστιονς ον Παττερν Αναλψσις ανδ Μασηινε Ιντελλιγενςε* 23.3 (2001), σσ. 257–267. ΔΟΙ: 10.1109/34.910878.
- [3] Α. Ζελνικ-Μανορ και Μ. Ιρανι. ‘Έεντ-βασεδ αναλψσις οφ ιδεο’. Στο: *Προσεεδινγς οφ τηε 2001 IEEE δμπτυερ Σοσιετψ δνφερενςε ον δμπτυερ ίσιον ανδ Παττερν Ρεσογνιτιον. ΄΄ΠΡ 2001*. Τόμ. 2. 2001, σσ. Π–Π. ΔΟΙ: 10.1109/΄΄ΠΡ.2001.990935.
- [4] Λαπτε και Λινδεβεργ. ‘Σπασε-τιμε ιντερεστ ποιנטς’. Στο: *Προσεεδινγς Νιντη IEEE Ιντερνατιοναλ δνφερενςε ον δμπτυερ ίσιον*. 2003, 432–439 ολ.1. ΔΟΙ: 10.1109/Ι΄΄΄.2003.1238378.
- [5] Αλπερ Ψιλμαζ και Μυβαρακ Σηαη. ‘Αστιονς σκετςη: α νοελ αστιον ρεπρεσε-ντατιον’. Στο: *2005 IEEE δμπτυερ Σοσιετψ δνφερενςε ον δμπτυερ ίσιον ανδ Παττερν Ρεσογνιτιον (΄΄ΠΡ’05)*. Τόμ. 1. 2005, 984–989 ολ. 1. ΔΟΙ: 10.1109/΄΄ΠΡ.2005.58.
- [6] Μ. Βλανκ κ.ά. ‘Αστιονς ας σπασε-τιμε σηαπες’. Στο: *Τεντη IEEE Ιντερνατιοναλ δνφερενςε ον δμπτυερ ίσιον (Γ΄΄΄’05) δλυμε 1*. Τόμ. 2. 2005, 1395–1402 έλ. 2. ΔΟΙ: 10.1109/Ι΄΄΄.2005.28.
- [7] Π. Δολλαφ κ.ά. ‘Βεηαιορ ρεσογνιτιον ια σπαρσε σπατιο-τεμποραλ φεατυρες’. Στο: *2005 IEEE Ιντερνατιοναλ Ωορκσηοπ ον ίσιναλ Συρειλλαλνσε ανδ Περφορμανςε Εαλυατιον οφ Τραςκινγ ανδ Συρειλλαλνσε*. 2005, σσ. 65–72. ΔΟΙ: 10.1109/΄΄ΣΠΕΤΣ.2005.1570899.
- [8] Ε. Σηεσητμαν και Μ. Ιρανι. ‘Σπασε-τιμε βεηαιορ βασεδ ζορρελατιον’. Στο: *2005 IEEE δμπτυερ Σοσιετψ δνφερενςε ον δμπτυερ ίσιον ανδ Παττερν Ρεσογνιτιον (΄΄ΠΡ’05)*. Τόμ. 1. 2005, 405–412 ολ. 1. ΔΟΙ: 10.1109/΄΄ΠΡ.2005.328.
- [9] Ψ. Σηεικη, Μ. Σηεικη και Μ. Σηαη. ‘Εξπλορινγ τηε σπασε οφ α ημμαν αστιον’. Στο: *Τεντη IEEE Ιντερνατιοναλ δνφερενςε ον δμπτυερ ίσιον (Γ΄΄΄’05) δλυμε 1*. Τόμ. 1. 2005, 144–149 έλ. 1. ΔΟΙ: 10.1109/Ι΄΄΄.2005.90.

- [10] Α. Ψιλμαζ και Μ. Σηαη. ‘Ρεσογνιζινγ ηυμαν αςτιονς ιν ιδεος αςχειρεδ βψ υνςαλιβρατεδ μοινγ ςαμερας’. Στο: *Τεντη IEEE Ιντερνατιοναλ δνφερενςε ον δμπττερ ιςιον (Γ^{ωω}05) δλυμε 1*. Τόμ. 1. 2005, 150–157 δλ. 1. doi: 10.1109/Γ^{ωω}.2005.201.
- [11] Θυαν άρλος Νιεβλες, Ηονγςηενγ Ωανγ και Φει Φει Λι. ‘Υνςυπερισεδ Λε-αρνινγ οφ Ηυμαν Αςτιον άτεγοριες Υςινγ Σπατιαλ-Τεμποραλ Ωορδς.’ Στο: τόμ. 79. Σεπτ. 2006, ςς. 1249–1258.
- [12] Μ. Σ. Ρψοο και Θ. Κ. Αγγαρωαλ. ‘Σεμαντις Υνδερςτανδινγ οφ δντινυεδ ανδ Ρεςυρςιε Ηυμαν Αςτιυτιες’. Στο: *18τη Ιντερνατιοναλ δνφερενςε ον Παττερν Ρεςογνιτιον (ΓΠΡ06)*. Τόμ. 1. 2006, ςς. 379–378. doi: 10.1109/Γ^{ωω}ΠΡ.2006.1043.
- [13] Ψ. Κε, Ρ. Συκτηανκαρ και Μ. Ηεβερετ. ‘Σπατιο-τεμποραλ Σηαπε ανδ Φλωω δρρελατιον φορ Αςτιον Ρεςογνιτιον’. Στο: *2007 IEEE δνφερενςε ον δμπττερ ιςιον ανδ Παττερν Ρεςογνιτιον*. 2007, ςς. 1–8. doi: 10.1109/Γ^{ωω}ΠΡ.2007.383512.
- [14] Μ. Δ. Ροδριγεζ, Θ. Αημεδ και Μ. Σηαη. ‘Αςτιον ΜΑ^ωΗ α ςπατιο-τεμποραλ Μαξιμου Αεραγε δρρελατιον Ηειγητ φιλτερ φορ αςτιον ρεσογνιτιον’. Στο: *2008 IEEE δνφερενςε ον δμπττερ ιςιον ανδ Παττερν Ρεςογνιτιον*. 2008, ςς. 1–8. doi: 10.1109/Γ^{ωω}ΠΡ.2008.4587727.
- [15] Π. Φ. Φελςενςζωαλβ κ.ά. ‘Οβθεςτ Δετεςτιον ωιτη Διςκριμινατιελψ Τραινεδ Παρτ-Βασεδ Μοδελς’. Στο: *IEEE Τρανςαςτιονς ον Παττερν Αναλψςις ανδ Μαςηινε Ιντελλιγενςε 32.9* (2010), ςς. 1627–1645. doi: 10.1109/ΤΠΑΜΙ.2009.167.
- [16] Θ.Κ. Αγγαρωαλ και Μ.Σ. Ρψοο. ‘Ηυμαν Αςτιυτψ Αναλψςις: Α Ρειεω’. Στο: *Α^ωΜ δμπτ. Συρ. 43.3* (Απρ. 2011), 16:1–16:43. ΙΣΣΝ: 0360-0300. doi: 10.1145/1922649.1922653. ΤΡΑ: ηττπ://doi.αςμ.οργ/10.1145/1922649.1922653.
- [17] Σ. Θι κ.ά. ‘3Δ δνολυτιοναλ Νευραλ Νετωορκς φορ Ηυμαν Αςτιον Ρεσογνιτιον’. Στο: *IEEE Τρανςαςτιονς ον Παττερν Αναλψςις ανδ Μαςηινε Ιντελλιγενςε 35.1* (2013), ςς. 221–231. doi: 10.1109/ΤΠΑΜΙ.2012.59.
- [18] Σ. Μανεν, Μ. Γυλλαυμιν και Λ. “. Γοολ. ‘Πριμε Οβθεςτ Προποςαλς ωιτη Ρανδομιζεδ Πριμς Αλγοριτημ’. Στο: *2013 IEEE Ιντερνατιοναλ δνφερενςε ον δμπττερ ιςιον*. 2013, ςς. 2536–2543. doi: 10.1109/Γ^{ωω}.2013.315.
- [19] Ψ. Τιαν, Ρ. Συκτηανκαρ και Μ. Σηαη. ‘Σπατιοτεμποραλ Δεφορμαβλε Παρτ Μοδελς φορ Αςτιον Δετεςτιον’. Στο: *2013 IEEE δνφερενςε ον δμπττερ ιςιον ανδ Παττερν Ρεςογνιτιον*. 2013, ςς. 2642–2649. doi: 10.1109/Γ^{ωω}ΠΡ.2013.341.
- [20] Θ.Ρ.Ρ. Υιθλινγς κ.ά. ‘Σελεςτιε Σεαρςη φορ Οβθεςτ Ρεσογνιτιον’. Στο: *Ιντερνατιοναλ Θουρναλ οφ δμπττερ ιςιον* (2013). doi: 10.1007/ς11263-013-0620-5. ΤΡΑ: ηττπ://ωω.ηυηπελεν.νλ/πυβλιςατιονς/σελεςτιεΣεαρςηΔραφτ.πδφ.

- [21] Ρ. Γιρσηκς κ.ά. 'Ριξη Φεατυρε Ηιεραρςηιε ϑορ Αςςυρατε Οβθεςτ Δετεςτιον ανδ Σεμαντις Σεγμεντατιον'. Στο: *2014 IEEE δνφερενςε ον δμπτυερ ίςιον ανδ Παττερν Ρεςογνιτιον*. 2014, ϑς. 580–587. ΔΟΙ: 10.1109/~~PP.2014.81.
- [22] Μ. Θαιν κ.ά. 'Αςτιον Λοςαλιζατιον ωιτη Τυβελετς φρομ Μοτιον'. Στο: *2014 IEEE δνφερενςε ον δμπτυερ ίςιον ανδ Παττερν Ρεςογνιτιον*. 2014, ϑς. 740–747. ΔΟΙ: 10.1109/~~PP.2014.100.
- [23] Α. Καρπατηψ κ.ά. 'Λαργε-Σςαλε ίδεο "λαςςιφιςατιον ωιτη δνολυτιοναλ Νευραλ Νετωορκς'. Στο: *2014 IEEE δνφερενςε ον δμπτυερ ίςιον ανδ Παττερν Ρεςογνιτιον*. 2014, ϑς. 1725–1732. ΔΟΙ: 10.1109/~~PP.2014.223.
- [24] Δαν Ονεατα κ.ά. 'Σπατιο-Τεμποραλ Οβθεςτ Δετεςτιον Προποςαλς'. Στο: Σεπτ. 2014. ΔΟΙ: 10.1007/978-3-319-10578-9~48.
- [25] Καρεν Σιμονψαν και Ανδρεω Ζιςσερμαν. 'Τωο-ςτρεαμ ϑονολυτιοναλ νετωορκς ϑορ αςτιον ρεςογνιτιον ιν ιδεος'. Στο: *Αδανςες ιν Νευραλ Ινφορματιον Προςεςσινγ Σψςτεμς*. 2014, ϑς. 568–576.
- [26] Ω. ηεν και Θ. Θ. δρςο. 'Αςτιον Δετεςτιον βψ Ιμπλιςιτ Ιντεντιοναλ Μοτιον "λυςτερινγ'. Στο: *2015 IEEE Ιντερνατιοναλ δνφερενςε ον δμπτυερ ίςιον (I^{ωω})*. 2015, ϑς. 3298–3306. ΔΟΙ: 10.1109/I~~.2015.377.
- [27] Θαν ". αν Γεμερτ κ.ά. 'ΑΠΤ: Αςτιον λοςαλιζατιον προποςαλς φρομ δενςε τραθεςτοριες'. Στο: *Προςεεδινγς οφ τηε Βριτιση Μαςηινε ίςιον δνφερενςε (BM^{ωω})*. BM^{ωω} Α Πρεςς, 2015, ϑς. 177.1–177.12. ISBN: 1-901725-53-7. ΔΟΙ: 10.5244/~.29.177. ΓΡΛ: ηττπς://δξ.δοι.οργ/10.5244/~.29.177.
- [28] Γ. Γκιοςζαρι και Θ. Μαλικ. 'Φινδινγ αςτιον τυβες'. Στο: *2015 IEEE δνφερενςε ον δμπτυερ ίςιον ανδ Παττερν Ρεςογνιτιον (~~PP)*. 2015, ϑς. 759–768. ΔΟΙ: 10.1109/~~PP.2015.7298676.
- [29] Θοε Ψυε-Ηει Νγ κ.ά. 'Βεψονδ ϑηορτ ϑνιππετς: Δεεπ νετωορκς ϑορ ιδεο ϑλαςςιφιςατιον'. Στο: *2015 IEEE δνφερενςε ον δμπτυερ ίςιον ανδ Παττερν Ρεςογνιτιον (~~PP)*. 2015, ϑς. 4694–4702. ΔΟΙ: 10.1109/~~PP.2015.7299101.
- [30] Κ. Σοομρο, Η. Ιδρεες και Μ. Σηαη. 'Αςτιον Λοςαλιζατιον ιν ίδεος τηρουγη δντεζτ Ωαλκ'. Στο: *2015 IEEE Ιντερνατιοναλ δνφερενςε ον δμπτυερ ίςιον (I^{ωω})*. 2015, ϑς. 3280–3288. ΔΟΙ: 10.1109/I~~.2015.375.
- [31] Δ. Τραν κ.ά. 'Λεαρνινγ Σπατιοτεμποραλ Φεατυρες ωιτη 3Δ δνολυτιοναλ Νετωορκς'. Στο: *2015 IEEE Ιντερνατιοναλ δνφερενςε ον δμπτυερ ίςιον (I^{ωω})*. 2015, ϑς. 4489–4497. ΔΟΙ: 10.1109/I~~.2015.510.
- [32] Π. Ωεινζαεπφελ, Ζ. Ηαρςηαουι και ". Σςημιδ. 'Λεαρνινγ το Τραςκ ϑορ Σπατιο-Τεμποραλ Αςτιον Λοςαλιζατιον'. Στο: *2015 IEEE Ιντερνατιοναλ δνφερενςε ον δμπτυερ ίςιον (I^{ωω})*. 2015, ϑς. 3164–3172. ΔΟΙ: 10.1109/I~~.2015.362.
- [33] Γ. Ψυ και Θ. Ψυαν. 'Φαςτ αςτιον προποςαλς ϑορ ηυμαν αςτιον δετεςτιον ανδ ϑεαρςη'. Στο: *2015 IEEE δνφερενςε ον δμπτυερ ίςιον ανδ Παττερν Ρεςογνιτιον*. 2015, ϑς. 1302–1311. ΔΟΙ: 10.1109/~~PP.2015.7298735.

- [34] Ρ. Γ. Ίνβις, Θ. Έρβεεκ και Ξ. Σζημιδ. ‘Ωεακλψ Συπερισεδ Οβθεστ Λοσαλιζατιον ωιτη Μυλτι-Φολδ Μυλτιπλε Ινστανζε Λεαρνινγ’. Στο: *IEEE Τρανσαστιονς ον Παττερν Αναλψσις ανδ Μασηινε Ιντελλιγεन्ςε* 39.1 (2016), σσ. 189–203. ΔΟΙ: 10.1109/ΤΠΑΜΙ.2016.2535231.
- [35] Ξ. Φειζητενηοφερ, Α. Πινζ και Α. Ζισσερμαν. ‘δνολυτιοναλ Τωο-Στρεαμ Νετωορκ Φυσιον φορ Ίδεο Αςτιον Ρεζογνιτιον’. Στο: *2016 IEEE δνφερενςε ον δμρυτερ Ίσιον ανδ Παττερν Ρεζογνιτιον (ΨΠΡ)*. 2016, σσ. 1933–1941. ΔΟΙ: 10.1109/ΨΨΠΡ.2016.213.
- [36] Κ. Ηε κ.ά. ‘Δεεπ Ρεσιδυαλ Λεαρνινγ φορ Ιμαγε Ρεζογνιτιον’. Στο: *2016 IEEE δνφερενςε ον δμρυτερ Ίσιον ανδ Παττερν Ρεζογνιτιον (ΨΠΡ)*. 2016, σσ. 770–778. ΔΟΙ: 10.1109/ΨΨΠΡ.2016.90.
- [37] Πασσαλ Μεττες, Θαν Ξ. αν Γεμερτ και έες Γ. Μ. Σνοεκ. ‘Σποτ Ον: Αςτιον Λοσαλιζατιον φορμ Ποιντλψ-Συπερισεδ Προποσαλς’. Στο: *δΡΡ αβς/1604.07602* (2016). αρΞι: 1604.07602. ΥΡΑ: ηττπ://αρξι.οργ/αβς/1604.07602.
- [38] Ξιαοθιανγ Πενγ και δρδελια Σζημιδ. ‘Μυλτι-ρεγιον τωο-στρεαμ Ρ-ΨΝΝ φορ αςτιον δετεςτιον’. Στο: *ΕΨΨ-Ευροπεαν δνφερενςε ον δμρυτερ Ίσιον*. Τόμ. 9908. Λεςτυρε Νοτες ιν δμρυτερ Σςιενςε. Αμστερδαμ, Νετηρελανδς: Σπρινγκερ, Οκτ. 2016, σσ. 744–759. ΔΟΙ: 10.1007/978-3-319-46493-0ΨΨ45. ΥΡΑ: ηττπς://ηαλ.ινρι.α.φορ/ηαλ-01349107.
- [39] Συμαν Σαηα κ.ά. ‘Δεεπ Λεαρνινγ φορ Δετεςτινγ Μυλτιπλε Σπαζε-Τιμε Αςτιον Τυβες ιν Ίδεοσ’. Στο: *δΡΡ αβς/1608.01529* (2016). αρΞι: 1608.01529. ΥΡΑ: ηττπ://αρξι.οργ/αβς/1608.01529.
- [40] Λιμιν Ωανγ κ.ά. ‘Τεμποραλ Σεγμεντ Νετωορκς: Τοωαρδς Γοοδ Πραςτιςες φορ Δεεπ Αςτιον Ρεζογνιτιον’. Στο: *δΡΡ αβς/1608.00859* (2016). αρΞι: 1608.00859. ΥΡΑ: ηττπ://αρξι.οργ/αβς/1608.00859.
- [41] Πηλιππε Ωεινζαεπφελ, Ξιαερ Μαρτιν και δρδελια Σζημιδ. ‘Τοωαρδς Ωεακλψ-Συπερισεδ Αςτιον Λοσαλιζατιον’. Στο: *δΡΡ αβς/1605.05197* (2016). αρΞι: 1605.05197. ΥΡΑ: ηττπ://αρξι.οργ/αβς/1605.05197.
- [42] Βοωεν Ζηανγ κ.ά. ‘Ρεαλ-τιμε Αςτιον Ρεζογνιτιον ωιτη Ενηανςεδ Μοτιον έςτορ ΨΝΝς’. Στο: *δΡΡ αβς/1604.07669* (2016). αρΞι: 1604.07669. ΥΡΑ: ηττπ://αρξι.οργ/αβς/1604.07669.
- [43] Ηαρχιρατ Σ. Βεηλ κ.ά. ‘Ινςρεμενταλ Τυβε δνστρυςτιον φορ Ηυμαν Αςτιον Δετεςτιον’. Στο: *δΡΡ αβς/1704.01358* (2017). αρΞι: 1704.01358. ΥΡΑ: ηττπ://αρξι.οργ/αβς/1704.01358.
- [44] Θ. άρρειρα και Α. Ζισσερμαν. ‘Χυο άδις, Αςτιον Ρεζογνιτιον; Α Νεω Μοδελ ανδ τηε Κινετις Δατασετ’. Στο: *2017 IEEE δνφερενςε ον δμρυτερ Ίσιον ανδ Παττερν Ρεζογνιτιον (ΨΠΡ)*. 2017, σσ. 4724–4733. ΔΟΙ: 10.1109/ΨΨΠΡ.2017.502.
- [45] Αλι Διβα κ.ά. ‘Τεμποραλ 3Δ δνΝετς: Νεω Αρςηιτεςτυρε ανδ Τρανσφερ Λεαρνινγ φορ Ίδεο Ψλαςσιφισατιον’. Στο: *δΡΡ αβς/1711.08200* (2017). αρΞι: 1711.08200. ΥΡΑ: ηττπ://αρξι.οργ/αβς/1711.08200.

- [46] Θ. Δοναηγε κ.ά. ‘Λονγ-Τερμ Ρεσυρρεντ δνολυτιοναλ Νετωορκς φορ Ίσυαλ Ρεζογνιτιον ανδ Δεσκριπτιον’. Στο: *IEEE Τρανσαςτιονς ον Παττερν Αναλυσις ανδ Μασηινε Ιντελλιγενςε* 39.4 (2017), σσ. 677–691. Δοι: 10.1109/ΤΠΑΜΙ.2016.2599174.
- [47] Ροηιτ Γιρδηαρ και Δεα Ραμαναν. ‘Αττεντιοναλ Ποολινγ φορ Αςτιον Ρεζογνιτιον’. Στο: *δΡΡ* αβς/1711.01467 (2017). αρΞι: 1711.01467. ΤΡΑ: ηττπ://αρξι.οργ/αβς/1711.01467.
- [48] Κ. Ηαφα, Η. Καταοκα και Ψ. Σατοη. ‘Λεαρνινγ Σπατιο-Τεμποραλ Φεατυρες ωιτη 3Δ Ρεσιδυαλ Νετωορκς φορ Αςτιον Ρεζογνιτιον’. Στο: *2017 IEEE Ιντερνατιοναλ δνφερενςε ον δμπτυερ Ίσιον Ωορκσηοψς (Γ^{ωω}Ω)*. 2017, σσ. 3154–3160. Δοι: 10.1109/Γ^{ωω}Ω.2017.373.
- [49] Γ. Ηυανγ κ.ά. ‘Δενσελψ δννεστεδ δνολυτιοναλ Νετωορκς’. Στο: *2017 IEEE δνφερενςε ον δμπτυερ Ίσιον ανδ Παττερν Ρεζογνιτιον (ΨΡΡ)*. 2017, σσ. 2261–2269. Δοι: 10.1109/ΨΡΡ.2017.243.
- [50] ηιη-Ψαο Μα κ.ά. ‘ΤΣ-ΛΣΤΜ ανδ Τεμποραλ-Ινζεπτιον: Εξπλοιτινγ Σπατιο-τεμποραλ Δψναμικς φορ Αςτινψ Ρεζογνιτιον’. Στο: *δΡΡ* αβς/1703.10667 (2017). αρΞι: 1703.10667. ΤΡΑ: ηττπ://αρξι.οργ/αβς/1703.10667.
- [51] Π. Μεττες και Ξ. Γ. Μ. Σνοεκ. ‘Σπατιαλ-Αωαρε Οβθεστ Εμβεδδινγς φορ Ζερο-Σηοτ Λογαλιζατιον ανδ Ήλασσιφικατιον οφ Αςτιονς’. Στο: *2017 IEEE Ιντερνατιοναλ δνφερενςε ον δμπτυερ Ίσιον (Γ^{ωω})*. 2017, σσ. 4453–4462. Δοι: 10.1109/Γ^{ωω}.2017.476.
- [52] Σ. Ρεν κ.ά. ‘Φαστερ Ρ-ΨΝΝ: Τοωαρδς Ρεαλ-Τιμε Οβθεστ Δετεκτιον ωιτη Ρεγιον Προποσαλ Νετωορκς’. Στο: *IEEE Τρανσαςτιονς ον Παττερν Αναλυσις ανδ Μασηινε Ιντελλιγενςε* 39.6 (2017), σσ. 1137–1149. Δοι: 10.1109/ΤΠΑΜΙ.2016.2577031.
- [53] Κ. Σοομρο και Μ. Σηαη. ‘Υνσυπερισεδ Αςτιον Δισοερψ ανδ Λογαλιζατιον ιν Ίδεος’. Στο: *2017 IEEE Ιντερνατιοναλ δνφερενςε ον δμπτυερ Ίσιον (Γ^{ωω})*. 2017, σσ. 696–705. Δοι: 10.1109/Γ^{ωω}.2017.82.
- [54] Ψι Ζηυ κ.ά. ‘Ηιδδεν Τωο-Στρεαμ δνολυτιοναλ Νετωορκς φορ Αςτιον Ρεζογνιτιον’. Στο: *δΡΡ* αβς/1704.00389 (2017). αρΞι: 1704.00389. ΤΡΑ: ηττπ://αρξι.οργ/αβς/1704.00389.
- [55] Αλι Διβα κ.ά. ‘Τεμποραλ 3Δ δνΝετς Υσινγ Τεμποραλ Τρανσιτιον Λαψερ’. Στο: *2018 IEEE δνφερενςε ον δμπτυερ Ίσιον ανδ Παττερν Ρεζογνιτιον Ωορκσηοψς, ΨΡΡ Ωορκσηοψς 2018, Σαλτ Λακε Ήτψ, ΥΤ, ΥΣΑ, Θυνε 18-22, 2018*. 2018, σσ. 1117–1121. ΤΡΑ: ηττπ://οπεναςεσς.τηεσφ.ομ/ζοντεντ^Ψ ^Ψςπρ^Ψ ^Ψ2018^Ψ ^Ψορκσηοψς/ω19/ητμλ/Διβα^Ψ ^ΨΤεμποραλ^Ψ ^Ψ3Δ^Ψ ^ΨδνΝετς^Ψ ^ΨΨΨΡ^Ψ ^Ψ2018^Ψ ^Ψπαπερ.ητμλ.
- [56] Ροηιτ Γιρδηαρ κ.ά. ‘Α Βεττερ Βασελινε φορ ΑΨΑ’. Στο: *δΡΡ* αβς/1807.10066 (2018). αρΞι: 1807.10066. ΤΡΑ: ηττπ://αρξι.οργ/αβς/1807.10066.
- [57] Ξ. Γυ κ.ά. ‘ΑΨΑ: Α Ίδεο Δατασετ οφ Σπατιο-Τεμποραλλψ Λογαλιζεδ Ατομικς Ίσυαλ Αςτιονς’. Στο: *2018 IEEE/ΨΦ δνφερενςε ον δμπτυερ Ίσιον ανδ Παττερν Ρεζογνιτιον*. 2018, σσ. 6047–6056. Δοι: 10.1109/ΨΨΡΡ.2018.00633.

- [58] Μιςηελλε Γυο κ.ά. ‘Νευραλ Γραπη Ματσηηγ Νετωορκς φορ Φεωσηοτ 3Δ Αςτιον Ρεςογνιτιον’. Στο: *Τηε Ευροπεαν δνφερενςε ον δμπτυερ ίσιον (E^{αω})*. 2018.
- [59] Θ. Ηε κ.ά. ‘Γενερις Τυβελετ Προποσαλς φορ Αςτιον Λοςαλιζατιον’. Στο: *2018 IEEE Ωντερ δνφερενςε ον Αππλιςατιονς οφ δμπτυερ ίσιον (ΩΑ^{αω})*. 2018, σς. 343–351. Δοι: 10.1109/ΩΑ^{αω}.2018.00044.
- [60] Ψυ Κογγ και Ψυν Φυ. ‘Ηυμαν Αςτιον Ρεςογνιτιον ανδ Πρεδιστιον: Α Συρεψ’. Στο: *δΡΡ αβς/1806.11230* (2018). αρΞι: 1806.11230. ΓΡΑ: ηττπ://αρξι.οργ/αβς/1806.11230.
- [61] Δ. Τραν κ.ά. ‘Α Ξλοσερ Λοοκ ατ Σπατιοτεμποραλ δνολυτιονς φορ Αςτιον Ρεςογνιτιον’. Στο: *2018 IEEE/ΞΦ δνφερενςε ον δμπτυερ ίσιον ανδ Παττερν Ρεςογνιτιον*. 2018, σς. 6450–6459. Δοι: 10.1109/ΞΞΠΡ.2018.00675.