



(12) CERERE DE BREVET DE INVENȚIE

(21) Nr. cerere: a 2022 00334

(22) Data de depozit: 15/06/2022

(41) Data publicării cererii:  
29/12/2023 BOPI nr. 12/2023

(71) Solicitant:  
• ARTIFICIAL INTELLIGENCE EXPERT  
S.R.L., STR.ALEXANDRU VLAHUȚĂ,  
BLOC LAMA C, NR.45, CLUJ-NAPOCA, CJ,  
RO

(72) Inventatori:  
• FLOARES ALEXANDRU,  
STR. ALEXANDRU VLAHUȚĂ, BL. LAMA C,  
AP. 45, CLUJ-NAPOCA, CJ, RO;  
• ZETY ADRIAN, STR.AUREL VLAICU,  
BL.17, AP.11, CLUJ-NAPOCA, CJ, RO

(54) SISTEM INTELIGENT DE ANALIZĂ AUTOMATĂ A DATELOR  
MICRORNA NGS

(57) Rezumat:

Invenția se referă la un sistem și la un algoritm de analiză a datelor microRNA brute, obținute prin tehnologia NGS (Next Generation Sequencing), cu ajutorul unui flux de lucru ce combină etape bioinformatică și de inteligență artificială pentru a obține biomarkeri relevanți pentru o problemă investigată și modele predictive pentru a o soluționa, cu aplicare în domenii precum zootehnia, agricultura, criminalistica și medicina. Algoritmul de analiză, conform invenției, cuprinde următoarele etape:

- preprocesare și controlul calității datelor, în care se verifică calitatea secvențelor și se produce un raport detaliat cu toate informațiile din conținutul unei probe analizate, iar raportul este analizat și se decide asupra filtrării secvențelor de calitate nesatisfăcătoare,
- îndepărtarea secvențelor adaptor,
- alinierea secvențelor pe o referință,
- cuantificarea miRNA, în care un algoritm va colapsa toate secvențele identice și le va număra, după care le va introduce într-o matrice care va fi completată ulterior cu numărul secvențelor din fiecare probă, obținând un profil complet de expresie al miRNA care este utilizat ca intrare într-un Expert AI al cărui flux de lucru cuprinde: analiza exploratorie și curățarea datelor, preprocesarea datelor înainte de antrenarea modelelor, dezvoltarea de modele predictive și optimizarea parametrilor algoritmilor de modelare.

Revendicări: 1  
Figuri: 2

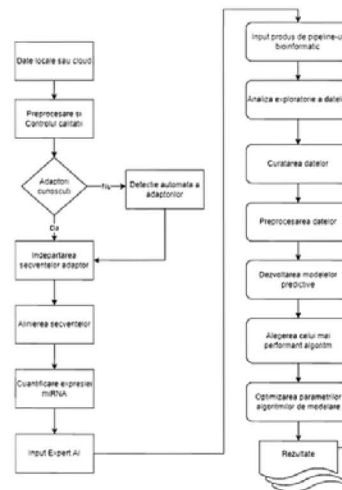


Fig. 2

Cu începere de la data publicării cererii de brevet, cererea asigură, în mod provizoriu, solicitantului, protecția conferită potrivit dispozițiilor art.32 din Legea nr.64/1991, cu excepția cazurilor în care cererea de brevet de invenție a fost respinsă, retrasă sau considerată ca fiind retrasă. Întinderea protecției conferite de cererea de brevet de invenție este determinată de revendicările conținute în cererea publicată în conformitate cu art.23 alin.(1) - (3).



## Descrierea

### Sistem Inteligent de Analiza Automata a Datelor microRNA NGS

Invenția se referă la un sistem și un algoritm de analiza a datelor microRNA brute, obținute prin tehnologiile Next Generation Sequencing (NGS), cu ajutorul unui flux de lucru ce combina etape Bioinformatică și de Inteligență Artificială (AI), pentru a obține biomarkeri relevanți pentru problema investigată și modele predictive performante, pentru a o soluționa.

Dat fiind rolul important al miRNA de reglare posttranslatională a activității genelor în întreaga lume vie a pluricelulelor, descoperirea biomarkerilor și dezvoltarea de modele predictive este importantă în domenii precum zootehnia, agricultura, criminalistica, și medicina.

MicroRNA (miRNA) sunt molecule mici de ARN necodificatoare de proteine care joacă un rol regulator important în translația genelor prin degradarea sau blocarea unor molecule de ARN mesager. Acestea influențează numeroase procese biologice majore, incluzând diferențierea, proliferare, apoptoza, inflamația și metabolismul<sup>1</sup>. Datorită rolului proeminent pe care miRNA îl joacă în expresia genică și funcționalitatea normală a organismelor, nu e surprinzător faptul că expresia lor aberantă poate duce la o multitudine de boli incluzând cancerul, bolile neurodegenerative, diabetul, condițiile cardiace, disfuncții ale rinichilor și ficatului, alterări ale sistemului imun. În plus, pe lângă contribuția la cauza a numeroase boli, microRNA pot fi utilizate și în terapiile țintite și în generarea de biomarkeri pentru diagnosticarea precoce a multor afecțiuni. Dintre acestea, un rol important îl joacă moleculele de miRNA circulante care pot prelua rolul de biomarkeri<sup>2</sup> fiind detectabile din probe de sânge, stand la baza metodelor neinvazive de detecție a numeroase boli<sup>3</sup>.

În oricare din domeniile de aplicabilitate, problema fundamentală, ce este departe de-a fi satisfăcător rezolvată, este aceeași. Utilizarea statisticii convenționale conduce la o listă de miRNA ce au o expresie diferită în două sau mai multe situații de interes distincte (expresie diferențială), de exemplu Cancer sau Normal. O astfel de listă este de foarte mică utilitate practică. Inteligența Artificială (utilizată în sistemul nostru) face posibilă obținerea unor modele predictive, cu valoare practică (de exemplu teste de diagnostic al cancerului cu acuratețe >95%). Simpla aplicare a AI,

<sup>1</sup> O'Brien J, Hayder H, Zayed Y, Peng C. Overview of microRNA biogenesis, mechanisms of actions, and circulation. *Front. Endocrinol.* 9, 402 (2018).

<sup>2</sup> Fabbri M (2010) miRNAs as molecular biomarkers of cancer. *Expert Rev* 10(4):435–444

<sup>3</sup> Chen X et al. (2008) Characterization of microRNAs in serum: a novel class of biomarkers for diagnosis of cancer and other diseases. *Cell Res* 18(10):997–1006

chiar a celor mai avansate tehnici, nu conduce automat la astfel de performanțe întrucât sistemele vii au o proprietate aparte - redundanță funcțională. Acesta este de regulă neglijată (nu și de către sistemul nostru) conducând la modele ce își degradează performanța când sunt aplicate la cazuri noi, diferite de cele folosite la învățare. Există numeroase domenii ce pot beneficia de pe urma analizei miRNA precum zootehnia, agricultura, medicina.

Asadar, printre principalele avantaje ale sistemului nostru se enumera și faptul că iau în considerare redundanță funcțională conducând la modele robuste ce generalizează bine la noi cazuri.

Datele brute miRNA nu vorbesc de la sine. Pentru a deveni instrumente practice utile, de exemplu un test de diagnostic al cancerului, etc., este necesară descoperirea din date a unor modele predictive cât mai performante (de exemplu, cu acuratețe >95%). Acest lucru este posibil doar prin aplicarea unor tehnici avansate de Inteligența Artificială (AI). Totuși, sistemele vii (plante, animale, ființe umane) prezintă o caracteristică aparte care face ca simpla aplicare a AI să nu dea rezultate satisfăcătoare - modele predictive robuste care generalizează bine (acuratețea nu se degradează) când sunt aplicate la cazuri noi, diferite de cele utilizate la învățare. Sistemele vii au o înaltă redundanță funcțională. Simplu spus, aceeași funcție se poate realiza la nivel molecular în mai multe moduri. Aceasta oferă sistemelor vii rezistență la o multitudine de alterări posibile. Cel mai adesea, redundanța funcțională este neglijată și se încearcă obținerea unui subset de miRNA cât mai mic.

În prezent, sistemele, algoritmi și fluxurile de lucru pentru descoperirea de biomarkeri miRNA, implementate în pachete software comerciale - Qiagen CLC Workbench, Illumina miRNAs Analysis, Partek Flow, - sau open source - Chipster, COMPSRA, Prost!, sau DIANA-mAP - se finalizează prin analiza expresiei diferențiale și nu prin modele predictive performante. Se obține doar o listă de miRNA exprimate diferit între cele două sau mai multe situații de interes (de exemplu, Cancer și Normal). Această listă este apoi trunchiată arbitrar și aplicată cazurilor noi. Se compară expresia miRNA a cazurilor noi cu această listă pentru încadrarea în una din clasele investigate. Prin contrast, soluția noastră descoperă modele predictive. Expresia miRNA a noilor cazuri este introdusă în acestea și se obține încadrare în una din clase. Deși rețeta de analiză ce provine de la datele miRNA NGS brute este similară la toate soluțiile existente, fiind oarecum standard, soluția noastră introduce un întreg flux de analiză cu IA ce conduce la soluții de înalt pragmatism și performanță.

Toate sistemele conventionale existente in prezent prezinta urmatoarele probleme:

- Metoda statistica de obtinere a subsetului de miRNA relevanti pentru problema data, din întregul set de miRNA cunoscuți și determinati cantitativ de echipamentul NGS si de pipeline-ul bioinformatic, nu primește feedback de la acuratețea cu care este rezolvata problema (de exemplu, diagnosticarea cancerului). Aceasta are impact asupra acuratetei soluției.
- Conține un număr foarte mare de miRNA exprimate diferit. Determinarea cantitativa a tuturor acestora, pentru aplicarea la situații concrete, este atat costisitoare cat și inutilă.
- Pentru reducerea costurilor se alege un prag arbitrar de expresie diferentia pentru a scurta lista.
- Expresia biomarkerilor miRNA, obtinuti prin trunchierea listei complete trebuie apoi comparata, mai mult sau mai puțin manual, cu expresia acelorași miRNA la un nou caz la care se aplica.
- Drept rezultat, soluțiile tehnice cunoscute au adesea acuratețe scăzută, sunt costisitoare și dificil de aplicat.

Deși etapele de procesare bioinformatica a datelor miRNA NGS sunt similare la toate soluțiile (inclusiv a noastră), fiind standard, invenția noastră înlătura dezavantajele soluțiilor tehnice cunoscute, înlocuind analiza expresiei diferențiale (statistica conventionala) cu un întreg flux de lucru bazat pe Inteligența Artificială (AI). Mai mult, pentru a putea fi folosită cu ușurință și de către cei fără cunoștințe de programare sau AI, face apel la metode de automatizare a Machine Learning (AutoML). Am numit acest flux automat de aplicare a Inteligenței Artificiale Expert AI pentru ai sugera functionalitatea.

Toate acestea fac posibilă obținerea unor biomarkeri relevanti și a unor modele predictive cu valoare practica (de exemplu teste de diagnostic al cancerului cu acuratețe >95%).

În toate domeniile de aplicabilitate, problema importantă dar dificila este de a descoperi biomarkeri și modele predictive cu o înaltă acuratețe, simpla analiza a expresiei diferențiale fiind total nesatisfăcătoare. De fapt, tehnicile de AI din sistemul nostru vor elimina oricum variabilele ce nu diferă cantitativ între situațiile de discriminat intrucat acestea nu sunt informative pentru predicție. Aplicand AI, selectarea variabilelor se face cu ajutorul unui feedback de la acuratețea predicției. In plus exista o paleta larga de algoritmi AI ceea ce permite experimentarea și alegerea celui mai bun spre deosebire de analiza statistica conventionala. Aceste avantaje sunt posibile datorită înlocuirii

metodei statistice de analiza a expresiei diferențiale cu un flux de lucru bazat pe Inteligenței Artificiale Automate (AutoML).

Pentru descoperirea biomarkerilor miRNA relevanti și transformarea acestora în instrumente pragmatice, de soluționare a problemelor în domeniile menționate mai sus, prin intermediul unor modele predictive, introducem metodologii de Inteligența Artificială adaptate specificului sistemelor vii, caracterizate prin redundanța funcțională la nivel molecular (miRNA aici), în locul analizei statistice convenționale - expresia diferențială a miRNA. În plus, prin implementarea Inteligenței Artificiale Automate, invenția poate fi utilizată cu ușurință și de către persoane fără cunoștințe avansate de programare.

Se dă, în continuare, un exemplu de realizare a invenției, în legătură cu figurile 1, 2 și 3, care reprezintă:

- Figura 1. Sistem - vedere de ansamblu
- Figura 2. Workflow Bioinformatic specific datelor NGS și Expert AI

Se dă, în continuare un exemplu de realizare a invenției pe un set de date publice miRNA-Seq NGS de cancer la san preluate din The Cancer Genome Atlas (<https://portal.gdc.cancer.gov/projects/TCGA-BRCA>) cu o cohorta de 1079 de pacienți cu cancer la san și un număr egal de cazuri non-cancer. Trebuie subliniat foarte clar ca pașii necesari obținerii datelor - colectarea sangelui/ plasmei, extracția ARN, conversia ARN în librării miRNA pentru NGS și secvențierea probelor - NU fac obiectul invenției. Invenția noastră se aplică DOAR datelor nu și probelor pacienților. Etapele realizării invenției sunt următoarele:

**Preprocesare și controlul calității datelor.** Se verifică calitatea secvențelor și se produce un raport detaliat cu toate informațiile despre conținutul probei. Rapoartele pentru probe individuale sunt ulterior incluse într-un raport general ce conține informații despre întregul experiment prezentând detalii precum conținutul de baze azotate (Adenina, Guanina, Citozina, Uracil), prezenta și structura secvențelor adaptor (în situația în care acestea nu au fost îndepărtate anterior, proces ce este oferit de către unele centre de secvențiere). Raportul este analizat și se decide asupra filtrării secvențelor de calitate nesatisfăcătoare.

**Îndepărtarea secvențelor adaptor.** Acestea pot fi cunoscute anterior, facilitățile de secvențiere oferind de obicei o documentație în care menționează secvențele adaptor utilizate, sau pot fi identificate prin utilizarea unor algoritmi specializați. Deodată cu eliminarea secvențelor adaptor se poate introduce și un prag de lungime pentru secvențele ce vor fi obținute. Dat fiind că secvențele

miRNA au o lungime medie de 22 de nucleotide, se va seta un prag cu o valoare între 22-30 nucleotide pentru a păstra doar moleculele de interes și pentru a scurta timpul de procesare în etapa de aliniere (se va utiliza o valoare mai mare de 22 pentru a acomoda eventualele erori rămase din etapa de filtrare).

**Aliniere a secvențelor pe o referință.** Se va utiliza cea mai recentă variantă publicată și adnotată a genomului uman. Secvențele preprocesate vor fi aliniate pe genom utilizând cei mai performanți algoritmi de aliniere după care vor fi adnotate utilizând cea mai recentă versiune a bazelor de date pentru miRNA (ex: miRBase).

**Cuantificarea miRNA.** Un algoritm va colapsa toate secvențele identice și le va număra după care le va introduce într-o matrice (tabel). Această matrice va fi completată ulterior cu numărul secvențelor din fiecare probă obținând astfel un tabel cu profilul complet de expresie al miRNA. Profilul de expresie al miRNA e utilizat ca input în Expertul AI.

**Analiza Exploratorie și Curățarea Datelor.** Primul pas al workflow-ului expertului AI este analiza exploratorie a datelor. Aceasta începe cu verificarea formei datelor în care se observă numărul de linii și coloane. Una din aceste coloane va deveni ținta pentru care va fi realizată predicția utilizând algoritmi AI. Datele de input și output trebuie să aibă formatul corect - coloanele miRNA conțin numere reale, iar ținta - diagnosticul aici - este variabila categorială.

Se vor identifica eventualele date lipsă și procentul acestora, eventualele anomalii prezente în date și se va proceda la o analiză statistică sumară (distribuția datelor, valori minime, maxime, medii, etc.). Ulterior se vor verifica corelațiile dintre variabilele prezente, astfel se poate observa potențialul ca unele variabile să fie mai importante în analiză și construcția modelelor predictive.

#### **Preprocesare datelor înainte de antrenarea modelelor:**

Identificare valorilor lipsă și imputarea acestora cu metode avansate.

Identificarea de Outliers prin multiple metode și îndepărtarea acestora.

Identificare de eventuale denumiri multiple pentru aceeași coloană și corectare.

Logaritmarea în baza 2 a datelor. Se consideră că distribuția datelor se apropie de cea Gaussiană, prin această transformare.

Normalizarea. Se alege una din următoarele metode: scalarea datelor între 0 și 1, sau între -1 și 1 sau, mai frecvent, standardizarea – media datelor devine 0 iar deviația standard devine 1.

### Dezvoltarea modelelor predictive

Prin utilizarea metodelor AutoML, modelele predictive vor fi dezvoltate automat. Totusi, vor fi cerute anumite setari din partea utilizatorului, facilitate de folosirea unui notebook (de exemplu, Jupyter Notebooks).

Se vor incarca librăriile și pachetele necesare mediului de lucru pentru Expertul AI (fluxul de lucru AI).

Se vor defini perechile input - valorile de expresie a miRNA si output - coloana tinta (diagnostic) cu valorile "Cancer" si "Non-Cancer").

Se definesc seturile de antrenare si testare. Daca numărul cazurilor permite, se definește si cel de al 3 lea set (de validare). Cel mai adesea, numărul cazurilor este relativ mic, motiv pentru care se definesc doar un set de antrenare și testare și se utilizează cross validation pentru antrenare. De exemplu, se împart cazurile de antrenare în 10 subseturi, învățarea are loc pe 9 dintre ele iar testarea pe al 10 lea, substerile fiind permutate.

In acest exemplu, performanta clasificarii binare Cancer - Non-Cancer poate fi reprezentată de mai multor metrici precum Acuratețea, AUC (Area Under the Curve), Precizia, Kappa si MCC (coeficient de corelatie Matthews).

Se antrenează toate modele disponibile in librariile incarcate si se evaluează performanța acestora, rezultand o ierarhie de modele.

### Optimizarea parametrilor algoritmilor de modelare:

Orice algoritm de modelare are el însuși o serie de parametri (de exemplu, adancimea arborilor decizionali pentru ansamblele de arbori decizionali). Valorile obișnuite (default) ale acestora pot sa nu conducă la performanta maxima posibila, dar pot fi optimizate. În acest scop se utilizează subsetul de validare care, cand sunt putine cazuri, poate fi un subset al setului de antrenare. Daca performanța obținută este superioară se retin valorile corespunzătoare ale parametrilor algoritmului.

Prin aplicarea invenției se obțin următoarele avantaje:

- Se descoperă biomarkeri cu adevărat relevanți, intrucat acestia sunt doar o componenta a unor modelelor predictive a caror performanta este masurata. Daca modelul bazat pe un set de biomarkeri are o acuratete de peste 95%, in mod clar markerii sunt relevanti. Daca, inasa, acuratețea este doar de 60% de exemplu, valoare biomarkerilor descoperiți este cvasi nula, intrucat alegerea clasei la întâmplare are o acuratețe de 50%, destul de apropiata.

- Se descoperă modele predictive de mare performanță, întrucât acestea extrag toată informația conținută în inputuri, relativ la outputuri, fiind și direct aplicabile la noi cazuri, spre deosebire de simplele liste.
- Îmbunătățiri semnificative ale eficienței însoțite de scăderea costurilor în domeniile de aplicabilitate.



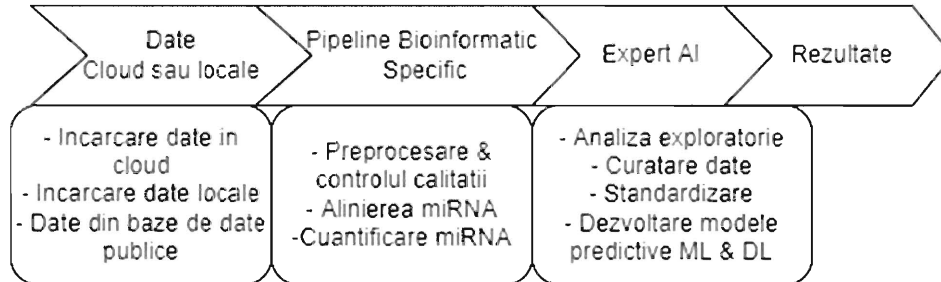
## Revendicari

Inventia permite descoperirea de biomarkeri și modele predictive cu o înaltă acuratețe. Tehnicile AI din inventie elimina variabilele ce nu diferă cantitativ între situațiile de discriminat păstrand doar variabilele informative pentru predicție, selectarea acestora se desfășoară utilizând un mecanism de feedback de la acuratețea predicției. În plus prin utilizarea abordărilor de tip AutoML se testează în mod automat numeroși algoritmi AI fapt care permite experimentarea și alegerea celui mai bun algoritm față de abordările de analiză statistică convențională.

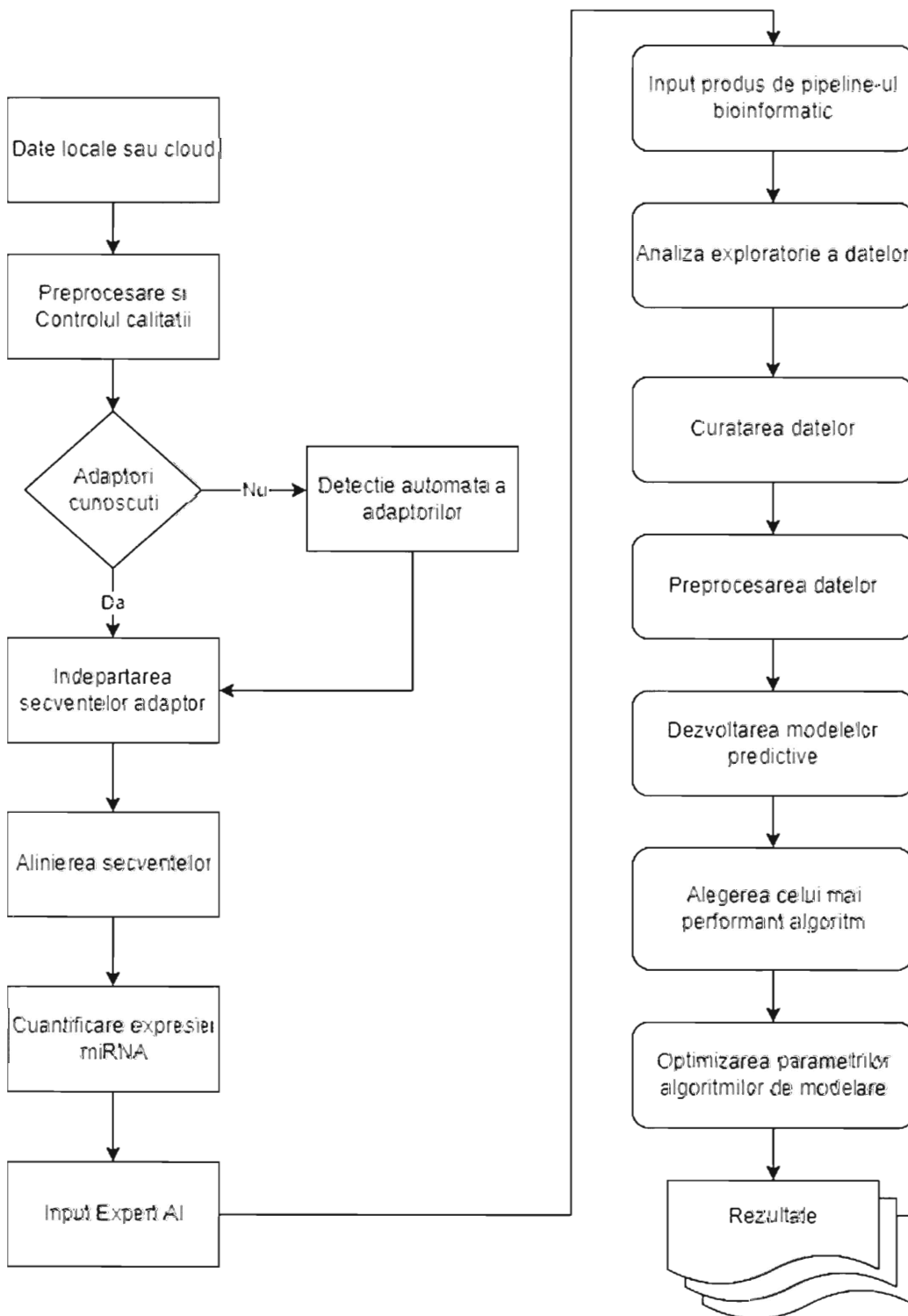
Invenția noastră constă într-un “Sistem Inteligent de Analiză a Datelor miRNA NGS” capabil de a prelua datele brute miRNA produse de echipamentele de laborator NGS, de a le preprocesa din punct de vedere bioinformatic (în modul standard) și de a descoperi biomarkeri în scopul soluționării diverselor probleme din domeniile de aplicabilitate (de exemplu, diagnosticul cancerului), caracterizată prin aceea că:

1. Metodelor statistice de analiză a expresiei diferențiale sunt înlocuite cu un flux de lucru bazat pe Inteligența Artificială iar
  2. Acest flux de lucru (Expert AI) este automatizat,
- ceea ce facilitează dramatic descoperirea biomarkerilor relevanți și a unor modele predictive cu înaltă acuratețe, oferind soluții pragmatice problemelor complexe din domeniile de aplicabilitate.

## Desene



**Figura 1.** Sistem - vedere de ansamblu



**Figura 2.** Workflow Bioinformatic specific datelor NGS si Expert AI