# NIP search and analysis in Metazoans

Jörg Lehmann

Bioinformatics Group
Department of Computer Science
University of Leipzig

25th TBI Winterseminar in Bled, 2010

**Outline**

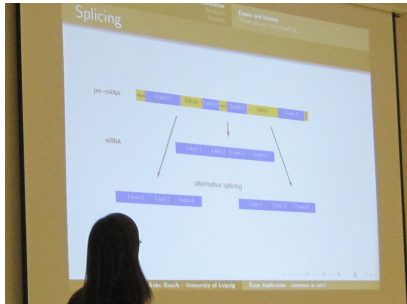**1 Motivation**
- The NIP Idea
- Previous Work

**2 NIPs in Metazoa**
- Data Compilation and NIP Extraction
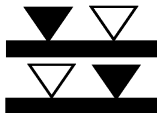- NIP Distribution Results

**Introns? Evolution?**
**A quick reminder.**

- spliceosomal intron positions often conserved across species
- intron positions as genome-level character



(source: Langenberger 2010)

**Outline**

# Near Intron Positions (NIPs)

IntronPos1                    IntronPos2

```
ref protein          _____119-0_____                    _____124-0_____
     Dsec GCTGATTCCCCTGCCATCAAGGAGAATCGCgtaaggttg//atttattcttgtagGCCTTTGGAGTGCAG          ACCATTTCCGGCACGGGTGCCCTGCGCATT
     Dsim GCTGATTCCCCTGCCATCAAGGAGAATCGCgtaaggttg//atttattcttgtagGCCTTTGGAGTGCAG          ACCATTTCGGGCACGGGTGCCCTGCGCATT
     Dmel GCCGATTCCCCAGCCATCAAGGAGAATCGCgtaaggttg//caattttcttgtagGCCTTTGGAGTGCAG          ACCATTTCGGGCACGGGTGCCCTGCGCGTT
     Dere GCTGATTCCCCAGCCATCAAGGAGAATCGCgtaaggcta//aatttttcttgtagGCTTTTGGAGTGCAG          ACCATTTCGGGCACAGGTGCCCTGCGTATT
     Dyak GCTGATTCCCCGGCCATCAAGGAGAATCGGgtaaggttg//cgatttttcttgtagGCCTTTGGAGTGCAG          ACCATTTCGGGCACGGGTGCCCTGCGTGTT
     Dana AGTGATTCTCCTGCTATTAAGGCCAATCGCgtaaggtag//tttttgtctggtagGCTTTTGGAGTTCAG          ACCATTTCGGGCACTGGTGCTCTGCGCATT
     Dper GCTGACTCCATAGCCATCAAGGAGAATCGCgtaatttta//cttcttgttcgcagGCTTTTGGGGTGCAG          ACTATATCCGGAACAGGTGCTCTGCGCGTG
     Dpse GCTGACTCCATAGCCATCAAGGAGAATCGCgtaatttta//cttcttgttcgcagGCTTTTGGGGTGCAG          ACTATATCCGGAACAGGTGCTCTGCGCGTG
     Dwil GCCGATTCCATAGCTCTTAAGGAGAAACGT                GCCTTTGGTGTCCAGgtgagttg//gatgcattttcagACTATATCTGGAACTGGTGCCTTGCGTGTG
     Dgri GCTGATTCCAATGCAATCAAAGAGAAACGAgttagtttt//tctgtcttctgcagGCTTTTGGTGTACAA          ACGATATCTGGAACAGGTGCAATTCGTGTC
     Dvir TCGGAATCAAAAGCATTAAAAGAGAAGCGTgtaagttt//ttattatttttgcagGCTTTTGGTGTGCAA          ACTATATCCGGCACAGGTGCGCTGCGTGTA
     Dmoj GCGGACTCAAATAGCATTAAAAGAGAAACGTgtaaggttg//tttttatattgcagGCCTTTGGTGTGCAA          ACGATTTCTGGAACAGGAGCTCTGCGCGTA
```

$<50nt$

The NIP Idea

# Phylogenetic Inference

Motivation
○○○●○○

NIPs in Metazoa
○○○○○○○○○○

Summary

Previous Work

**Outline**

## NIPs in Holometabolous Insects.

- NIP as phylogenetic marker (hypothesis testing)
- beetles more closely related to (butter-)flies than bees/wasps.



Krauss et al.: Near intron positions are reliable phylogenetic markers: an application to holometabolous insects. *Mol Biol Evol*, 25:821-830, 2008.

# Novel Introns in Drosophila.

- NIPs to identify recent intron gain and its mechanisms
- intron sliding and tandem duplication are causes for novel introns in *Drosophila*



(3) tandem exon duplication including a proto-splice site



(6) intron sliding



relocation of a preexisting intron within a gene
by movement of intron boundaries (mutations/indels) or
change of alternative splicing profiles

exonic DNA          intronic DNA

# Metazoan Species Overview.
## Supposed Cladogram.

**Outline**

**Compilation of Orthologous Genes.**

1. start: all *C. elegans* protein-coding genes with *D. melanogaster* 1:1 ortholog (`Ensembl/BioMart`) [**3331** genes]

2. retrieve `EnsemblCompara` orthologs [13 species]

3. `BLAST`-based ortholog annotation and CDS extraction [$> 28$ genomes]

4. exclusion of some species (too few/no protein data available)

## NIP Pipeline

1. Multiple codon alignment (Muscle, transAlign.pl, ...)
2. Extraction of alignment regions containing NIPs of distance < 32 nt (out: 45890 NIP regions from 3276 genes)
3. Filtering steps. Require:
   - amino-acid conservation score $\geq 0.75$
   - $\geq 6$ non-gap characters around introns
   - strict splice site consensus (g, ct, ag:act, a, g)

### Final Dataset

- **5616 NIPs** (from 2589 aln regions and 1220 genes)
- within **32 species**
- (1349 parsimony-informative NIPs)

Data Compilation and NIP Extraction

# Metazoan NIPs: A Nice Example Alignment.

```
(Dmel protein)                _____149-1_____       _____152-1_____
      Mbre ATGCCATGCTGGACGCCCCCGTATCTGCCGgttcgtag//tcgatgcgtggcagGTGTCGGAG                GCGCTGAAGCTGGCACTCTGACCTTTATGG
      Tadh GTACTTTTATGGACGCTCCAGTTTCTGGAGgtacaact//attcccgtaaatagGAGTTGTAG                CGGCTCGCGATGCGTTATTAACATTTATGG
      Hsap CAGTTTTCATGGATGCCCCCTGTTTCTGGTGgtaagtgg//cccaaatgattcagGTGTAGGAG               CTGCACGATCTGGGAACCTCACGTTTATGG
      Mdom CACTTTTCATGGATGCCCCTGTTTGGTGgtgagtaa//atcttttttctccagGGGTGGGAG               CTGCTCGAGCTGGGAACCTGACATTCATGG
     Mdom1 CACTTTTCCTGGATGTCCCTGTTTCTGGTG                        GG-------                -------CGTGGGAACCTGGCATTCATGG
      Drer CTGTTTTCATGGATGCCCCGGTGTCTGGAGgtgagatg//ttaatcttcctcagGTGTTGGTG                CTGCCACTTCTGGTAAGCTCACTTTTATGG
     Drer1 CAGTGTTCATGGACGCGCCGGTGTCAGGAGgtatgatg//tgtgttctgtctagGTGTTGGTG                CAGCCAGTTTGGCTAAACTCACTTTCTTGG
      Bflo GAGTCTACATGGACGCACCTGTTTCTGGGGgtatggta//tctacacattacagGTATTAAGG                 CAGCTGCTGCAGCCACACTAACCTTCATGG
      Csav CCAAGTATCTGGATGCACCAGTTTCTGGAGgtttgcat//attgttctacacagGTGTGGGTG                CCGCTCAAGCAGGAACTTTAACTTTCATGG
      Isca -------------------CTTTTACGGGGtatcacc//cttttatattctagGTGTGAATG                 CAGCCAAGGCTGGGACACTTACATTCATGG
      Dmel CTCGGTTTATCGATGCCCCCGTTTCCGGCG                  GAGTTCCCG                     GCGCCGAGCAGGCCACCCTCACCTTCATGG
      Agam CAACGTTCGTCGATGCGCCCGTGTCCGGCG                  GTGTGCCCG                     GTGCAAAGAACGCGACGCTCACGTTCATGG
      Aaeg CTACCTTCGTCGATGCCCCGGTTTCCGGTG                  GTGTTCCCG                     GAGCTCAGAATGCTACTCTGACCTTCATGG
      Cpip CGACCTTTGTCGATGCGCCCGTTTCCGGTG                  GAGTTCCCG                     GAGCTAAGAACGCCACGCTGACGTTCATGG
      Tcas TCAAGTTTCTCGACGCTCCGGTCTCAGGTG                  GAGTCACGG                     GGGCGGAGGCCGGGACCCTCACTTTTATGG
      Phum CTTCATTTGTTGATGCTCCTGTGTCTGGAGgtgagtta//cggtttttttaagGTGTCATTG                 GCGCAAAAGATGGTACTCTGACTTTTATGG
      Dpul TGATTTTCATTGATGCTCCAGTTTCTGGAGgtattaaa//acggttattgacagGTGTGATGG                CTGCTAAAGCTGGAACTTTGACTTTCATGG
      Lgig CAGTTTATCTTGATGCACCAGTCTCAGGGGgtaagatc//taatttctttccagGTGTAAATG                CAGCTAGAGATGCTTTACTCACATTCATGG
      Ccap CGGTCTTCATGGATGCCCCTGTGTCTGGAGgtaggccg//tttgcatcgtgcagGAGTGGTGG                CCGCTAGAGACGCCCTCTTGACCTTCATGG
      Hrob CCACGTTTATTGATGCCCCTGTTTCTGGAGgtaacata//tgattgaacaaaagGTGTGACTG                CAGCCAAGGAAGGTACTTTGACGTTTATGG
      Minc CTAGTTTTTTGGATGCCCTGTATCTGGGG                    GAGTTTTAGgttcctta//cttttaatttgtagGTGCAGAAAAGGCTACTTTAACGTTTATGA
      Cjap CCGAGTACATTGATGCGCCGTTTTCGGGAG                   GAGTTACAGgtgggaaa//ccttttttattccagGCGCCCAAAACGCCACACTAACATTCATGG
      Cele CCGAGTATATCGACGCCCCAATCTCCGGTG                   GAGTCACTGgtacgggga//tcataattttccagGCGCCCAGCAAGCTACACTAACATTTATGG
      Cbre CTGAATACATCGATGCTCCGATTTCCGGAG                   GGGTTACAGgtttgttg//acttaaaacttcagGAGCCCAACAAGCCACACTCACTTTTATGG
      Cbri CTGAATACATCGATTCTCCAATTTCCGGAG                   GCGTCACTG                     GAGCTCAACAAGCCACGTTGACATTTATGG
      Crem CAGAGTACATTGATGCCACCGATTTCTGGTG                  GTGTCACTGgtacggag//ccaacaatttccagGCGCCCAACAAGCGACTCTCACTTTCATGG
      Nvec CAACATATTTAGATGCACCTGTTTCGGGAGgtattttt//tgttctatctttagGGATTACAG                 CAGCAAAAGCAGGTACACTGACATTTATGG
      Hmag GTTCATATGTGGATGCACCTGTTTCAGGGGgtaaattt//tatcatatgtttagGTGTTAACG                 CTGCAAAAGAAGGAACATTGACAATTATGG
```

from ortholog set B0250.5 (Cele)
ref protein is Dmel (FBpp0085821/FBgn0034390)

Data Compilation and NIP Extraction

# Metazoan NIPs: A Nice Example Alignment.

## Pos0  <32nt  Pos1

```
(Dmel protein)                    _____149-1_____          _____152-1_____
      Mbre ATGCCATGCTGGACGCCCCCGTATCTGGCGgttcgtag//tcgatgcgtggcagGTGTCGGAG      GCGCTGAAGCTGGCACTCTGACCTTTATGG
      Tadh GTACTTTTATGGACGCTCCAGTTTCTGGAGgtacaact//attcccgtaaatagGAGTTGTAG       CGGCTCGCGATGCGTTATTAACATTTATGG
      Hsap CAGTTTTCATGGATGCCCCCTGTTTCTGGTGgtaagtgg//cccaaatgattcagGTGTAGGAG        CTGCACGATCTGGGAACCTCACGTTTATGG
      Mdom CACTTTTCATGGATGCCCCTGTTTCTGGTGgtgagtaa//atctttttctccagGGGTGGGAG         CTGCTCGAGCTGGGAACCTGACATTCATGG
     Mdom1 CACTTTTCCTGGATGTCCCTGTTTCTGGTG                    GG-------                 -------CGTGGGAACCTGGCATTCATGG
      Drer CTGTTTTCATGGATGCCCCGGTGTCTGGAGgtgagatg//ttaatcttcctcagGTGTTGGTG          CTGCCACTTCTGGTAAGCTCACTTTTATGG
     Drer1 CAGTGTTCATGGACGCGCCGGTGTCAGGAGgtatgatg//tgtgttctgtctagGTGTTGGTG          CAGCCAGTTGGCTAAACTCACTTTCTTGG
      Bflo GAGTCTACATGGACGCACCTGTGTCTGGGGgtatggta//tctacacattacagGTATTAAGG           CAGCTGCTGCAGCCACACTAACCTTCATGG
      Csav CCAAGTATCTGGATGCACCAGTTTCTGGAGgtttgcat//attgttctacacagGTGTGGGTG          CCGCTCAAGCAGGAACTTTAACTTTCATGG
      Isca --------------------CTTTTACGGGgtatcacc//cttttatattctagGTGTGAATG          CAGCCAAGGCTGGGACACTTACATTCATGG
      Dmel CTCGGTTTATCGATGCCCCCGTTTCCGGCG                    GAGTTCCCG                 GCGCCGAGCAGGCCACCCTCACCTTCATGG
      Agam CAACGTTCGTCGATGCCGCCCGTGTCCGGCG                   GTGTGCCCG                 GTGCAAAGAACGCGACGCTCACGTTCATGG
      Aaeg CTACCTTCGTCGATGCCCCGTTTCCGGTG                     GTGTTCCCG                 GAGCTCAGAATGCTACTCTGACCTTCATGG
      Cpip CGACCTTTGTCGATGCGCCCGTTTCCGGTG                    GAGTTCCCG                 GAGCTAAGAACGCCACGCTGACGTTCATGG
      Tcas TCAAGTTTCTCGACGCTCCGGTCTCAGGTG                    GAGTCACGG                 GGGCGGAGGCCGGGACCCTCACTTTTATGG
      Phum CTTCATTTGTTGATGCTCCTGTGTCTGGAGgtgagtta//cggtttttttaagGTGTCATTG            GCGCAAAAGATGGTACTCTGACTTTCATGG
      Dpul TGATTTTCATTGATGCTCCAGTTTCTGGAGgtattaaa//acggttattgacagGTGTGATGG           CTGCTAAAGCTGGAACTTTGACTTTCATGG
      Lgig CAGTTTATCTTGATGCACCAGTCTCAGGGGgtaagatc//taatttctttccagGTGTAAATG           CAGCTAGAGATGCTTTACTCACATTCATGG
      Ccap CGGTCTTCATGGATGCCCCTGTGTCTGGAGgtaggccg//tttgcatcgtgcagGAGTGGTGG           CCGCTAGAGACGCCCTCTTGACCTTCATGG
      Hrob CCACGTTTATTGATGCCCCTGTTTCTGGAGgtaacata//tgattgaacaaaagGTGTGACTG           CAGCCAAGGAAGGTACTTTGACGTTTATGG
      Minc CTAGTTTTTTGGATGCCCCTGTATCTGGGG                GAGTTTTAGgttcctta//cttttaatttgtagGTGCAGAAAAGCTACTTTAACGTTTATGA
      Cjap CCGAGTACATTGATGCGCCCAATTTCGGGAG               GAGTTACAGgtgggaaa//ccttttttattccagGCGCCCAAAACGCCACACTAACATTCATGG
      Cele CCGAGTATATCGACGCCCCAATCTCCGGTG                GAGTCACTGgtacggga//tcataatttccagGCGCCCAGCAAGCTACACTAACATTTATGG
      Cbre CTGAATACATCGATGCTCCGATTTCCGGAG                GGGTTACAGgtttgtta//acttaaaacttccagGAGCCCAACAAGCCACACTCACTTTTATGG
      Cbri CTGAATACATCGATTCTCCAATTTCCGGAG                GCGTCACTG                 GAGCTCAACAAGCCACGTTGACATTTATGG
      Crem CAGAGTACATTGATGCACCGATTTCTGGTG                GTGTCACTGgtacggag//ccaacaatttccagGCGCCCAACAAGCGACTCTCACTTTCATGG
      Nvec CAACATATTTAGATGCACCTGTTTCTGGGGgtattttt//tgttctatctttagGGATTACAG            CAGCAAAAGCAGGTACACTGACATTTATGG
      Hmag GTTCATATGTGGATGCACCTGTTTCAGGGGgtaaattt//tatcatatgtttagGTGTTAACG            CTGCAAAAGAAGGAACATTGACAATTATGG
```

from ortholog set B0250.5 (Cele)
ref protein is Dmel (FBpp0085821/FBgn0034390)

## Metazoan NIPs: A Nice Example Alignment.

# Pos0        Pos1

```
(Dmel protein)                                        _____149-1_____         _____152-1_____
        Mbre H_A_M_L_D_A_P_V_S_G__gttcgtag//tcgatgcgtggcagG__V__G_             G_A_E_A_G_T_L_T_F_M_
        Tadh S_T_F_M_D_A_P_V_S_G__gtacaact//attcccgtaaatagG__V__V_             A_A_R_D_A_L_L_T_F_M_
        Hsap A_V_F_M_D_A_P_V_S_G__gtaagtgg//cccaaatgattcagG__V__G_             A_A_R_S_G_N_L_T_F_M_
        Mdom A_L_F_M_D_A_P_V_S_G__gtgagtaa//atcttttttctccagG__V__G_            A_A_R_A_G_N_L_T_F_M_
        Mdom1 P_L_F_L_D_V_P_V_S_G__                          G__-__-_          -_-__-_R_G_N_L_A_F_M_
        Drer A_V_F_M_D_A_P_V_S_G__gtgagatg//ttaatcttcctcagG__V_G_              A_A_T_S_G_K_L_T_F_M_
        Drer1 A_V_F_M_D_A_P_V_S_G__gtatgatg//tgtgttctgtctagG__V_G_             A_A_S_L_A_K_L_T_F_L_
        Bflo G_V_Y_M_D_A_P_V_S_G__gtatggta//tctacacattacagG__I_K_              A_A_A_A_A_T_L_T_F_M_
        Csav A_K_Y_L_D_A_P_V_S_G__gtttgcat//attgttctacacagG__V_G_              A_A_Q_A_G_T_L_T_F_M_
        Isca __-__-__-__-__-__-_L_L_R_gtatcacc//cttttatattctagG__V_N_          A_A_K_A_G_T_L_T_F_M_
        Dmel A_R_F_I_D_A_P_V_S_G__                          G_V_P_              G_A_E_Q_A_T_L_T_F_M_
        Agam A_T_F_V_D_A_P_V_S_G__                          G_V_P_              G_A_K_N_A_T_L_T_F_M_
        Aaeg A_T_F_V_D_A_P_V_S_G__                          G_V_P_              G_A_Q_N_A_T_L_T_F_M_
        Cpip A_T_F_V_D_A_P_V_S_G__                          G_V_P_              G_A_K_N_A_T_L_T_F_M_
        Tcas F_K_F_L_D_A_P_V_S_G__                          G_V_T_              G_A_E_A_G_T_L_T_F_M_
        Phum T_S_F_V_D_A_P_V_S_G__gtgagtta//cggtttttttaagG__V_I_               G_A_K_D_G_T_L_T_F_M_
        Dpul M_I_F_I_D_A_P_V_S_G__gtattaaa//acggttattgacagG__V_M_              A_A_A_A_G_T_L_T_F_M_
        Lgig S_V_Y_L_D_A_P_V_S_G__gtaagatc//taatttctttccagG__V_N_              A_A_R_D_A_L_L_T_F_M_
        Ccap S_V_F_M_D_A_P_V_S_G__gtaggccg//tttgcatcgtgcagG__V_V_              A_A_R_D_A_L_L_T_F_M_
        Hrob S_T_F_I_D_A_P_V_S_G__gtaacata//tgattgaacaaaagG__V_T_              A_A_K_E_G_T_L_T_F_M_
        Minc A_S_F_L_D_A_P_V_S_G__                      G_V_L__gttcctta//cttttaatttgtagG_A_E_K_A_T_L_T_F_M_
        Cjap A_E_Y_I_D_A_P_I_S_G__                      G_V_T__gtgggaaa//ccttttattccagG_A_Q_N_A_T_L_T_F_M_
        Cele A_E_Y_I_D_A_P_I_S_G__                      G_V_T__gtacggga//tcataatttccagG_A_Q_Q_A_T_L_T_F_M_
        Cbre A_E_Y_I_D_A_P_I_S_G__                      G_V_T__gtttgtta//acttaaaacttcagG_A_Q_Q_A_T_L_T_F_M_
        Cbri A_E_Y_I_D_S_P_I_S_G__                      G_V_T_              G_A_Q_Q_A_T_L_T_F_M_
        Crem A_E_Y_I_D_A_P_I_S_G__                      G_V_T__gtacggag//ccaacaatttccagG_A_Q_A_A_T_L_T_F_M_
        Nvec A_V_F_V_D_A_P_V_S_G__gtattttt//tgttctatcttttagG__I_T_            A_A_K_A_G_T_L_T_F_M_
        Hmag C_S_V_V_D_A_P_V_S_G__gtaaattt//tatcatatgtttagG__V_N_             A_A_K_E_G_T_L_T_F_I_M_
```

from ortholog set B0250.5 (Cele)
ref protein is Dmel (FBpp0085821/FBgn0034390)

Motivation
○○○○○○

NIPs in Metazoa
○○○○●○○○○○

Summary

Data Compilation and NIP Extraction

## Metazoan NIPs: Example Character Matrix.

```
'B0250.1.47-0_48-1'  'B0.
'B0250.5.149-1_152-1' 'B
MATRIX
Mbre    1????????
Tadh    ?00?0???

    . . .

Hrob    ????0???
Bmal    ?00??????
Cjap    ?00?1???
Cele    ????100?
Cbre    ?00??00?
Crem    ?00?100?
Cbri    ?00??00?
Sman    ?1?0????
Nvec    ?00?0?11
Hmag    ?00?0???
Mdom    000?0?11
Drer    000?0???
```
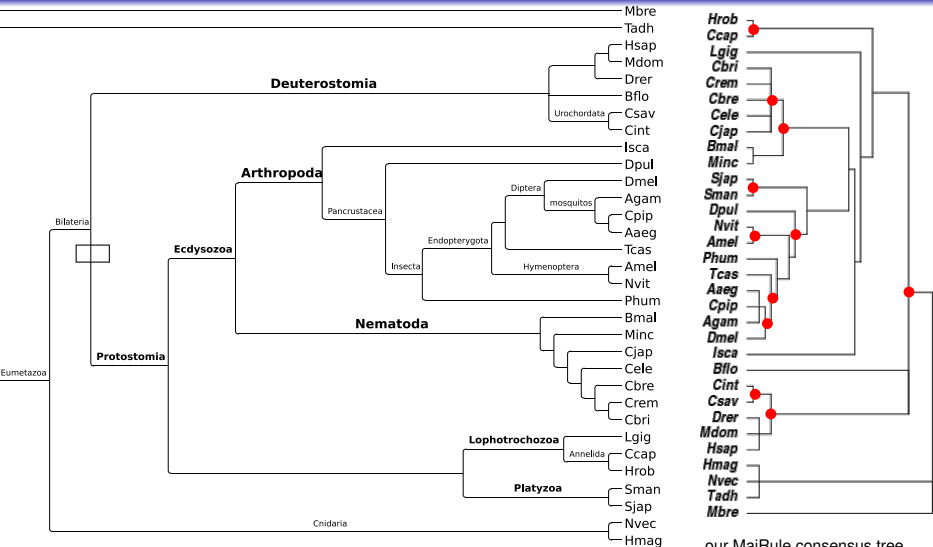
**Outline**

NIP Distribution Results

# Parsimony Tree Search with `PAUP*`.
## `hsearch addseq=random swap=tbr nreps=1000`
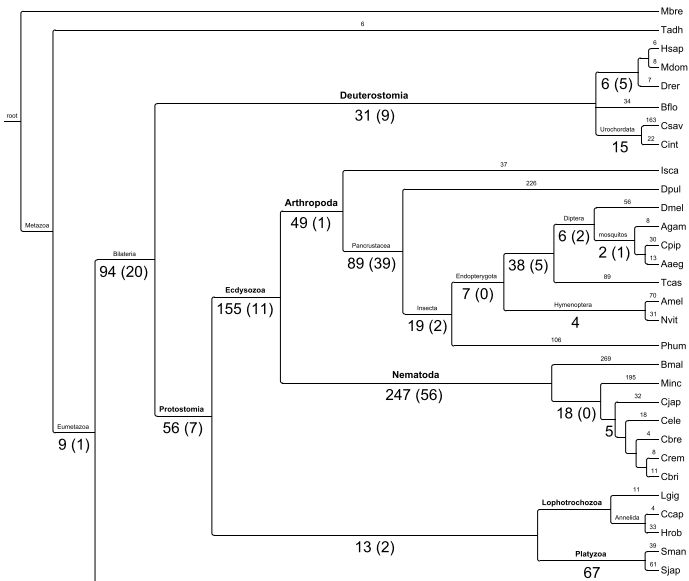


data: 1349 informative character          supposed tree: 1538 steps

our MajRule consensus tree
(17010 trees with 1510 steps)

# Synaptomorphic and Autapomorphic NIPs for Supposed Tree

**Open Questions.**

- position of Amphioxus (*Bflo*)
- position of the flatworms (Platyhelminthes, *Sjap*/*Sman*/*Smed*)

### Limitations with MP and NIP character

- **Long Branch Attraction** - solution: dense taxon sampling
- limited amount of informative character changes in some branches vs
- conflicting (but frequent) character changes leading to unresolved clades

**Where To Go Next?**

- improve/extend ortholog dataset
- evaluate specific parts of the Metazoan tree (testing phylogenetic hypotheses), e.g. within Caernorhabditis

## Summary

- NIP concept and application as phylogenetic marker and novel intron search
- NIP data for Metazoan seems promising

- To improve
  - more reliable and complete ortholog assignments (increase data and taxa)
  - extend automated NIP validation

## Thank You!



**Many thanks to:**

- Carina Eisenhardt (Genetics Group Uni Leipzig)
- Veiko Krauss
- Peter F. Stadler

## Thank you for your attention!

## Any questions?

**For Further Reading I**

Krauss V, Thümmler C, Georgi F, Lehmann J, Stadler PF, Eisenhardt C.
Near intron positions are reliable phylogenetic markers: an application to holometabolous insects.
*Mol Biol Evol*, 25:821-830, 2008.

Lehmann J, Eisenhardt C, Stadler PF, Krauss V.
Some novel intron positions in conserved Drosophila genes are caused by intron sliding or tandem duplication.
*BMC Evol Biol*, in revision.